

JASN

Plasma proteomics of renal function: a trans-ethnic meta-analysis and Mendelian randomization study

Journal:	<i>Journal of the American Society of Nephrology</i>
Manuscript ID	JASN-2020-07-1070.R1
Manuscript Type:	Original Article - Clinical Epidemiology
Date Submitted by the Author:	n/a
Complete List of Authors:	<p>Matías-García, Pamela; Helmholtz Zentrum München Research Unit Molecular Epidemiology; Helmholtz Zentrum München Institute of Epidemiology; Technical University of Munich; German Center for Cardiovascular Research (DZHK), partner site Munich Heart Alliance</p> <p>Wilson, Rory; Helmholtz Zentrum München Research Unit Molecular Epidemiology; Helmholtz Zentrum München Institute of Epidemiology</p> <p>Guo, Qi; University of Cambridge Department of Public Health and Primary Care, Cardiovascular Epidemiology Unit</p> <p>Zaghlool, Shaza; Weill Cornell Medicine - Qatar, Department of Physiology and Biophysics</p> <p>Eales, James; University of Manchester, Institute of Cardiovascular Sciences</p> <p>Xu, Xiaoguang; University of Manchester, Institute of Cardiovascular Sciences</p> <p>Charchar, Fadi; Federation University Australia, School of Health and Life Sciences; University of Leicester, Department of Cardiovascular Sciences; University of Melbourne, Department of Physiology</p> <p>Dormer, John; University Hospitals of Leicester, Cellular Pathology</p> <p>Schlosser, Pascal; Universitätsklinikum Freiburg, Institute of Genetic Epidemiology</p> <p>Elhadad, Mohamed; Helmholtz Zentrum München Research Unit Molecular Epidemiology; Helmholtz Zentrum München Institute of Epidemiology</p> <p>Nano, Jana; Helmholtz Zentrum München Institute of Epidemiology</p> <p>Sharma, Sapna; Helmholtz Zentrum München Research Unit Molecular Epidemiology; Helmholtz Zentrum München Institute of Epidemiology</p> <p>Peters, Annette; Helmholtz Zentrum München Institute of Epidemiology</p> <p>Fornoni, Alessia; Diabetes Research Institute, Miller School of Medicine, University of Miami, Medicine/Nephrology</p> <p>Mook-Kanamori, Dennis ; Leiden University Medical Center, Department of Clinical Epidemiology</p> <p>Winkelmann, Juliane; Helmholtz Zentrum München, Institute of Neurogenomics; Technical University of Munich, Department of Neurogenetics and Institute of Human Genetics</p> <p>Danesh, John; British Heart Foundation; University of Cambridge, Centre of Research Excellence</p> <p>Di Angelantonio, Emanuele ; University of Cambridge, Centre of Research Excellence ; University of Cambridge Department of Public Health and Primary Care, Cardiovascular Epidemiology Unit; National Institute for Health Research, Blood and Transplant Research Unit in</p>

	<p>Donor Health and Genomics Ouwehand, Willem; University of Cambridge, Department of Haematology; NHS Blood and Transplant, Cambridge Biomedical Campus; Wellcome Trust Sanger Institute Watkins, Nicholas; NHS Blood and Transplant, Cambridge Biomedical Campus Roberts, David; NHS Blood and Transplant, Oxford Centre; University of Oxford Radcliffe Department of Medicine Petrera, Agnese; Helmholtz Zentrum München Deutsches Forschungszentrum für Gesundheit und Umwelt, Research Unit Protein Science and Core Facility Proteomics Graumann, Johannes; Max Planck Institute for Heart and Lung Research, Scientific Service Group Biomolecular Mass Spectrometry; DZHK, partner site Rhine-Main Koenig, Wolfgang; German Center for Cardiovascular Research (DZHK); Deutsches Herzzentrum München; University of Ulm Institute for Epidemiology and Medical Biometrics Hveem, Kristian; Norwegian University of Science and Technology Faculty of Medicine and Health Sciences, K.G. Jebsen Centre for Genetic Epidemiology, Department of Public Health and Nursing, ; Norwegian University of Science and Technology Faculty of Medicine and Health Sciences, HUNT Research Centre Jonasson, Christian; Norwegian University of Science and Technology Faculty of Medicine and Health Sciences, K.G. Jebsen Centre for Genetic Epidemiology, Department of Public Health and Nursing, ; Norwegian University of Science and Technology Faculty of Medicine and Health Sciences, HUNT Research Centre Köttgen, Anna; Universitätsklinikum Freiburg, Institute of Genetic Epidemiology; Johns Hopkins Bloomberg School of Public Health, Department of Physiology and Biophysics Butterworth, Adam; University of Cambridge Department of Public Health and Primary Care, Cardiovascular Epidemiology Unit Prunotto, Marco; University of Geneva, School of Pharmaceutical Sciences Hauck, Stefanie; Helmholtz Zentrum München Deutsches Forschungszentrum für Gesundheit und Umwelt, Research Unit Protein Science and Core Facility Proteomics Suhre, Karsten; Weill Cornell Medicine - Qatar, Department of Physiology and Biophysics Gieger, Christian; Helmholtz Zentrum München Research Unit Molecular Epidemiology; Helmholtz Zentrum München Institute of Epidemiology; German Center for Cardiovascular Research (DZHK), partner site Munich Heart Alliance Tomaszewski, Maciej; University of Manchester, Institute of Cardiovascular Sciences Teumer, Alexander; Universitätsmedizin Greifswald, Institute for Community Medicine; German Center for Cardiovascular Research (DZHK), partner site Greifswald Waldenberger, Melanie; Helmholtz Zentrum München Research Unit Molecular Epidemiology; Helmholtz Zentrum München Institute of Epidemiology; German Center for Cardiovascular Research (DZHK), partner site Munich Heart Alliance</p>
Keywords:	Mendelian randomization, chronic kidney disease, glomerular filtration rate, Epidemiology and outcomes, plasma proteomics, biomarker discovery, testican-2, gene expression

SCHOLARONE™
Manuscripts

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

Authors: MatÃ-as-GarcÃ-a, Pamela; Wilson, Rory; Guo, Qi; Zaghlool, Shaza; Eales, James; Xu, Xiaoguang; Charchar, Fadi; Dormer, John; Schlosser, Pascal; Elhadad, Mohamed; Nano, Jana; Sharma, Sapna; Peters, Annette; Fornoni, Alessia; Mook-Kanamori, Dennis ; Winkelmann, Juliane; Danesh, John; Di Angelantonio, Emanuele ; Ouwehand, Willem; Watkins, Nicholas; Roberts, David; Petrera, Agnese; Graumann, Johannes; Koenig, Wolfgang; Hveem, Kristian; Jonasson, Christian; KÃttgen, Anna; Butterworth, Adam; Prunotto, Marco; Hauck, Stefanie; Suhre, Karsten; Gieger, Christian; Tomaszewski, Maciej; Teumer, Alexander; Waldenberger, Melanie

Title: Plasma proteomics of renal function: a trans-ethnic meta-analysis and Mendelian randomization study

Running title: Plasma proteomics of renal function

Manuscript Type: Original Article - Clinical Epidemiology

Funders: British Heart Foundation, (Grant / Award Number: 'PG/17/35/33001','PG/19/16/34270')
National Institute for Health Research, (Grant / Award Number:)
NIHR, (Grant / Award Number:)
Qatar National Research Fund (QNRF), (Grant / Award Number: 'NPRPC11-0115-180010')
Kidney Research UK, (Grant / Award Number: 'RP_013_20190305','RP_017_20180302')
Norwegian Institute of Public Health, (Grant / Award Number:)
Nord-TrÃndelag County Council, Central Norway Health Authority, (Grant / Award Number:)
NHS, (Grant / Award Number:)
Federal Ministry of Education and Research (BMBF), (Grant / Award Number:)
Qatar Foundation, (Grant / Award Number:)
Weill Cornell Medicine in Qatar, (Grant / Award Number:)
Faculty of Medicine, Norwegian University of Science and Technology, (Grant / Award Number:)
State of Bavaria, (Grant / Award Number:)

Financial Disclosure: No The authors declare the following competing interests: C.J. has received personal fees for research consultancy work from Pfizer and Bayer outside of the submitted work, and J.D. sits on the International Cardiovascular and Metabolic Advisory Board for Novartis (since 2010); the Steering Committee of UK Biobank (since 2011); the MRC International Advisory Group (ING) member, London (since 2013); the MRC High Throughput Science ~Omic Panel Member, London (since 2013); the Scientific Advisory Committee for Sanofi (since 2013); the International Cardiovascular and Metabolism Research and Development Portfolio Committee for Novartis; and the Astra Zeneca Genomics Advisory Board (2018). Other authors declare no competing interests.

Study Group/Organization Name: Human Kidney Tissue Studies

Study Group Members' Names: Wojciech Wystrychowski¹, Monika Szulinska², Andrzej Antczak³, Maciej Glyda⁴, Robert KrÃl¹, Joanna Zywiec⁵, Ewa Zukowska-Szzechowska⁶, Pawel Bogdanski²
¹ Department of General, Vascular and Transplant Surgery, Medical University of Silesia, Katowice, Poland
² Department of Treatment of Obesity, Metabolic Disorders and Clinical Dietetics, Poznan University of Medical Sciences, Poznan, Poland.

1
2
3 Department of Urology and Uro-oncology, Karol Marcinkowski University of Medical Sciences, Poznan,
4 Poland

5 4 Department of Transplantology and General Surgery Poznan, Collegium Medicum, Nicolaus Copernicus
6 University, Bydgoszcz, Poland

7 5 Department of Internal Medicine, Diabetology and Nephrology, Medical University of Silesia, Zabrze,
8 Poland

9 6 Department of Health Care, Silesian Medical College, Katowice, Poland

10
11
12 **Total number of words:** 4900

13
14 **Abstract:** Background. Studies on the relationship between renal function and the human plasma
15 proteome have identified several potential biomarkers. However, studies have been conducted largely
16 in European populations, and whether the associations between plasma proteins and kidney function
17 are causal has never been addressed.

18
19 **Methods.** A cross-sectional study of 993 plasma proteins and 2,882 participants of four studies of
20 European and admixed ancestries (KORA, INTERVAL, HUNT, QMDiab) was conducted to identify trans-
21 ethnic associations between eGFR/CKD and proteomic biomarkers. For the replicated associations, two-
22 sample bidirectional Mendelian randomization (MR) was used to investigate potential causal
23 relationships, followed by the analysis of gene expression in kidney.

24 **Results.** Fifty-seven plasma proteins were associated with eGFR, including two novel proteins, JAM-B
25 and contactin-4. Nineteen of these were additionally associated with CKD. The strongest inferred causal
26 effect was the positive effect of eGFR on testican-2, an effect in line with the known biological role of
27 this protein and the expression of the protein-coding gene of testican-2 (SPOCK2) in renal tissue. Finally,
28 we observed suggestive evidence of an effect of melanoma inhibitory activity (MIA), carbonic anhydrase
29 III and cystatin-M on eGFR.

30
31 **Conclusions.** In a discovery-replication setting, we identified 57 proteins trans-ethnically associated with
32 eGFR, including two potential novel biomarkers. Our findings with regard to causal relationships
33 represent an important stepping-stone in the establishment of testican-2 as a clinically relevant
34 physiological marker of kidney disease progression, and point to additional potential therapeutic targets
35 warranting further investigation.

36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

1
2
3 **Plasma proteomics of renal function: a trans-ethnic meta-analysis and Mendelian**
4 **randomization study**
5
6
7

8
9 **Significance statement**
10

11 Prior studies on the plasma proteome of renal function have identified several biomarkers, but
12 have lacked replication, were limited to European populations and/or did not investigate
13 causality. This paper describes, firstly, the identification of plasma proteomic biomarkers in the
14 largest cross-sectional study of renal function to date. Using four studies in a discovery-
15 replication setting, 57 protein biomarkers trans-ethnically associated with eGFR and/or CKD
16 were identified, two of which are novel. Investigations into causality using Mendelian
17 randomization provide suggestive evidence for a few proteins warranting further investigation as
18 therapeutic targets, and highlight testican-2 as a protein affected by renal function, an early
19 milestone in its establishment as a physiological marker of kidney disease progression with
20 potential clinical relevance.
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

Plasma proteomics of renal function: a trans-ethnic meta-analysis and Mendelian randomization study

Pamela R. Matías-García (ORCID iD: 0000-0002-3906-5405)^{1,2,3,4*}, Rory Wilson^{1,2}, Qi Guo⁵, Shaza B. Zaghlool⁶, [James M. Eales⁷](#), [Xiaoguang Xu⁷](#), [Fadi J. Charchar⁸⁻¹⁰](#), [John Dormer¹¹](#), Pascal Schlosser (0000-0002-8460-0462)^{7,12}, Mohamed A. Elhadad^{1,2,4}, Jana Nano^{2,8,13}, Sapna Sharma^{1,2}, Annette Peters^{2,4,8,13}, Alessia ~~Fornoni~~⁹[Fornoni¹⁴](#), Dennis O. Mook-[Kanamori¹⁰](#)[Kanamori¹⁵](#), Juliane Winkelmann^{11,12,16,17}, John Danesh^{13,14,15,16,17,18,23}, Emanuele Di Angelantonio^{13,14,15,16,17,18,22}, Willem H. Ouwehand^{14,19,20,21,19,24-26}, Nicholas A. [Watkins²⁰](#)[Watkins²⁵](#), David J. Roberts^{15,22,23,20,27,28}, [Agnese Petrera \(ORCID iD: 0000-0002-0583-6195\)²⁹](#), Johannes Graumann (ORCID iD: 0000-0002-3015-5850)^{24,25,30,31}, Wolfgang Koenig^{4,26,27,32,33}, Kristian Hveem^{28,29,34,35}, Christian Jonasson^{28,29,34,35}, Anna Köttgen (0000-0002-4671-3714)^{7,12,30,36}, Adam Butterworth⁵, Marco Prunotto (ORCID iD: 0000-0002-0203-0129)^{34,37}, [Stefanie M. Hauck \(ORCID iD: 0000-0002-1630-6827\)²⁹](#), Karsten Suhre⁶, Christian Gieger^{1,2,4}, [Maciej Tomaszewski^{7,38}](#), [Human Kidney Tissue Studies** \(names and affiliations provided after title page\)](#), Alexander Teumer^{32,33,39,40}, Melanie Waldenberger (ORCID iD: 0000-0003-0583-5093)^{1,2,4*}

¹ Research Unit of Molecular Epidemiology, Helmholtz Zentrum München, German Research Center for Environmental Health, 85764-Neuherberg, Bavaria, Germany.

² Institute of Epidemiology, Helmholtz Zentrum München, German Research Center for Environmental Health, Neuherberg, Bavaria, Germany.

³ TUM School of Medicine, Technical University of Munich, Munich, Germany.

⁴ German Center for Cardiovascular Research (DZHK), partner site Munich Heart Alliance, Munich, Germany.

⁵ Cardiovascular Epidemiology Unit, Department of Public Health and Primary Care, University of Cambridge, Strangeways Research Laboratory, Wort's Causeway, Cambridge, UK.

⁶ Department of Physiology and Biophysics, Weill Cornell Medicine-Qatar, Doha, Qatar.

1
2
3 [7 Division of Cardiovascular Sciences, University of Manchester, Manchester, UK.](#)

4
5 [8 School of Health and Life Sciences, Federation University Australia, Ballarat, VIC,](#)
6
7 [Australia.](#)

8
9 [9 Department of Cardiovascular Sciences, University of Leicester, Leicester, UK.](#)

10
11 [10 Department of Physiology, University of Melbourne, Melbourne, VIC, Australia.](#)

12
13 [11 University Hospitals of Leicester NHS Trust, Leicester, UK.](#)

14
15 [127](#) Institute of Genetic Epidemiology, Faculty of Medicine and Medical Center-University
16
17 of Freiburg, Freiburg, Germany.

18
19 [138](#) German Center for Diabetes Research (DZD), München-Neuherberg, Neuherberg,
20
21 Germany.

22
23 [149](#) Katz Family Division of Nephrology and Hypertension, Department of Medicine,
24
25 University of Miami Miller School of Medicine, Miami, FL, USA

26
27 [1540](#) Department of Clinical Epidemiology, Leiden University Medical Centre, Leiden, The
28
29 Netherlands.

30
31 [1644](#) Institute of Neurogenomics, Helmholtz Zentrum München, German Research Center
32
33 for Environmental Health, Neuherberg, Germany.

34
35 [1742](#) Department of Neurogenetics and Institute of Human Genetics, Technical University
36
37 of Munich, Munich, Germany.

38
39 [1843](#) British Heart Foundation Cardiovascular Epidemiology Unit, Department of Public
40
41 Health and Primary Care, University of Cambridge, Cambridge, UK

42
43 [1944](#) British Heart Foundation Centre of Research Excellence, University of Cambridge,
44
45 Cambridge, UK.

46
47 [2045](#) National Institute for Health Research Blood and Transplant Research Unit in Donor
48
49 Health and Genomics, University of Cambridge, Cambridge, UK.

50
51 [2146](#) National Institute for Health Research Cambridge Biomedical Research Centre,
52
53 University of Cambridge and Cambridge University Hospitals, Cambridge, UK.

54
55 [2247](#) Health Data Research UK Cambridge, Wellcome Genome Campus and University of
56
57 Cambridge, Cambridge, UK.

1
2
3 [2348](#) Department of Human Genetics, Wellcome Sanger Institute, Hinxton, UK.

4
5 [2449](#) Department of Haematology, University of Cambridge, Cambridge, UK.

6
7 [2520](#) NHS Blood and Transplant, Cambridge Biomedical Campus, Long Road, Cambridge,
8
9 UK.

10
11 [2624](#) Wellcome Sanger Institute, Hinxton, UK.

12
13 [2722](#) NHS Blood and Transplant-Oxford Centre, Level 2, John Radcliffe Hospital, Oxford,
14
15 UK.

16
17 [2823](#) Radcliffe Department of Medicine, University of Oxford, John Radcliffe Hospital,
18
19 Oxford, UK.

20
21 [29](#) Research Unit Protein Science and Core Facility Proteomics, Helmholtz Zentrum
22
23 München, German Research Center for Environmental Health, Neuherberg, Germany.

24
25 [3024](#) Scientific Service Group Biomolecular Mass Spectrometry, Max Planck Institute for
26
27 Heart and Lung Research, W.G. Kerckhoff Institute, Bad Nauheim, Germany.

28
29 [3125](#) German Centre for Cardiovascular Research (DZHK), partner site Rhine-Main, Max
30
31 Planck Institute of Heart and Lung Research, Bad Nauheim, Germany.

32
33 [3226](#) Deutsches Herzzentrum München, Technical University of Munich, Munich,
34
35 Germany.

36
37 [3327](#) Institute of Epidemiology and Medical Biometry, University of Ulm, Ulm, Germany.

38
39 [3428](#) K.G. Jebsen Centre for Genetic Epidemiology, Department of Public Health and
40
41 Nursing, Faculty of Medicine and Health Science, Norwegian University of Science and
42
43 Technology, Trondheim, Norway.

44
45 [3529](#) HUNT Research Centre, Faculty of Medicine, Norwegian University of Science and
46
47 Technology, Levanger, Norway.

48
49 [3630](#) Department of Epidemiology, Johns Hopkins Bloomberg School of Public Health,
50
51 Baltimore, MD, USA.

52
53 [3734](#) School of Pharmaceutical Sciences, University of Geneva, Geneva, Switzerland.

54
55 [38](#) Division of Medicine and Manchester Academic Health Science Centre, Manchester
56
57 University NHS Foundation Trust, Manchester, UK.

1
2
3 | [3932](#) Institute for Community Medicine, University Medicine Greifswald, Greifswald,
4
5 Germany.

6
7 | [4033](#) German Center for Cardiovascular Research (DZHK), partner site Greifswald,
8
9 Greifswald, Germany.

10
11
12
13 **Running title:** Plasma proteomics of renal function

14
15 **Abstract word count:** [231239](#)

16
17 **Main text word count:** [3,4354960](#)

18
19
20 **Corresponding author's information:** * Correspondence to Pamela R. Matías-García and
21
22 Melanie Waldenberger. Research Unit of Molecular Epidemiology - Institute of Epidemiology,
23
24 Helmholtz Zentrum München, Ingolstädter Landstraße 1, 85764 Munich/Neuherberg,
25
26 Germany; pamela.matias@helmholtz-muenchen.de, waldenberger@helmholtz-muenchen.de

27
28 **Keywords:** proteomics; Mendelian randomization; causal analysis; genetic variant;
29
30 biomarkers; chronic kidney disease; eGFR; testican-2; genetic epidemiology; plasma
31
32 proteins.
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

1
2
3 **Names and affiliations from authors in the group Human Kidney Tissue Studies****

4
5 **** The following investigators contributed to recruitment and/or phenotyping of human kidney**
6
7 **gene expression studies reported in ¹: Wojciech Wystrychowski¹, Monika Szulinska², Andrzej**
8 **Antczak³, Maciej Glyda⁴, Robert Król¹, Joanna Zywiec⁵, Ewa Zukowska-Szczechowska⁶,**
9 **Pawel Bogdanski²**

10
11
12
13 **1 Department of General, Vascular and Transplant Surgery, Medical University of Silesia,**
14 **Katowice, Poland**

15
16
17
18 **2 Department of Treatment of Obesity, Metabolic Disorders and Clinical Dietetics, Poznan**
19 **University of Medical Sciences, Poznan, Poland.**

20
21
22 **3 Department of Urology and Uro-oncology, Karol Marcinkowski University of Medical**
23 **Sciences, Poznan, Poland**

24
25
26
27 **4 Department of Transplantology and General Surgery Poznan, Collegium Medicum, Nicolaus**
28 **Copernicus University, Bydgoszcz, Poland**

29
30
31 **5 Department of Internal Medicine, Diabetology and Nephrology, Medical University of Silesia,**
32 **Zabrze, Poland**

33
34
35 **6 Department of Health Care, Silesian Medical College, Katowice, Poland**

Significance statement

Prior studies on the plasma proteome of renal function have identified several biomarkers, but have lacked replication, were limited to European populations and/or did not investigate causality. This paper describes, firstly, the identification of plasma proteomic biomarkers in the largest cross-sectional study of renal function to date. Using four studies ~~total~~ in a discovery-replication setting, 57 protein biomarkers trans-ethnically associated with eGFR and/or CKD were identified, two of which are novel. Investigations into causality using Mendelian randomization provide suggestive evidence for a few proteins warranting further investigation as therapeutic targets, and highlight testican-2 as a protein affected by renal function, an early milestone in its establishment as a physiological marker of kidney disease progression with potential clinical relevance.

Abstract

Background. Studies on the relationship between renal function and the human plasma proteome have identified several potential biomarkers. However, studies have been conducted largely in European populations, and whether the associations between plasma proteins and kidney function are causal has never been addressed.

Methods. A cross-sectional study of 993 plasma proteins and 2,882 participants of four studies of European and admixed ancestries (KORA, INTERVAL, HUNT, QMDiab) was conducted to identify trans-ethnic associations between eGFR/CKD and proteomic biomarkers. For the replicated associations, two-sample bidirectional Mendelian randomization (MR) was used to investigate potential causal relationships, followed by the analysis of gene expression in kidney.

Results. Fifty-seven plasma proteins were associated with eGFR, including two novel proteins, JAM-B and contactin-4. Nineteen of these were additionally associated with CKD. The strongest inferred causal effect ~~estimated by MR~~ was the positive effect of eGFR on testican-2, an effect in line with the known biological role of this protein and the glomeruli-specific expression of the protein-coding gene of testican-2 (*SPOCK2*) in renal tissue. Finally, we observed suggestive evidence of an effect of melanoma inhibitory activity (MIA), carbonic anhydrase III and cystatin-M on eGFR.

Conclusions. In a discovery-replication setting, we identified 57 proteins trans-ethnically associated with eGFR, including two potential novel biomarkers. Our findings with regard to causal relationships represent an important stepping-stone in the establishment of testican-2 as a clinically relevant physiological marker of kidney disease progression, and point to additional potential therapeutic targets warranting further investigation.

Introduction

The kidneys' ability to filter blood and maintain homeostasis is reflected in the glomerular filtration rate (GFR)². Serum creatinine, a filtration marker, is freely filtered by the glomerulus³ and its blood levels can be used to calculate estimated GFR (eGFR)⁴. Chronic kidney disease (CKD), characterized by reduced eGFR (<60 ml/min/m²) and proteinuria, has a global prevalence of 10% to 16%^{5, 6} and is expected to be increasingly common in aging populations². Increased serum creatinine is not evident until ~50% of the renal filtration function is lost⁷, making CKD a silent disease and creating a blind spot for early renal disease detection⁸. Its rising prevalence, in addition to the lack of therapeutic options⁸, imposes a significant burden on health systems and individuals worldwide^{2, 6}.

A number of biomarker research studies have been conducted in regard to early detection, diagnosis and/or progression prediction of kidney diseases^{7, 8}. Early efforts in renal function proteome research focused on urine biomarkers, though most studies were small and lacked replication^{7, 8}; combining multiple urinary biomarkers proved to be more successful (e.g. CKD273 classifier predicting deterioration and mortality)^{7, 9}. More recent studies have focused on blood, an easily accessible tissue that reflects the metabolic status of multiple organs. Its complex proteomic profile, however, requires sensitive and reliable techniques for its study. A promising tool in this regard is SOMAscan, a platform using DNA aptamers to measure hundreds of plasma proteomic biomarkers¹⁰. Although this platform has been successfully used in different epidemiological settings¹¹⁻¹⁵, renal disease has not been sufficiently investigated: prior studies have either tested a limited number of proteins¹⁶ or relied on small samples without replication^{10, 17}.

Moreover, prior studies on proteomic markers and renal function have not distinguished causality from correlation¹⁸. Mendelian randomization (MR), an instrumental variable analysis used to infer causal effects on a given outcome by relying on the random allocation of alleles at conception and genetic variation as a proxy to lifetime exposures, is an increasingly popular method used in genetic epidemiology studies to address causality^{19, 20}.

1
2
3 Here we present a cross-sectional study using a multiplexed aptamer-based
4 proteomics platform to investigate associations between 1095 plasma proteins and
5 GFR/CKD and other renal parameters in a discovery cohort (KORA, N=995), with replication
6 in three independent studies of European and admixed ancestry (INTERVAL, HUNT and
7 QMDiab; N =1,887). To better understand the biological significance of the identified
8 proteins, we conducted enrichment, protein and gene expression analyses across tissues, as
9 well as investigated their interconnection using protein-protein interaction (PPI) network
10 analysis. We also investigated causal effects between eGFR and the replicated proteins
11 using two-sample bidirectional Mendelian randomization, and further examined gene
12 expression in kidney tissue to further ~~explore~~ investigate the proteins ~~with the strongest~~
13 evidence of a causal effect.
14
15
16
17
18
19
20
21
22
23
24
25

26 **Methods**

27 ***Study populations***

28
29 The KORA study (Cooperative Health Research in the Region of Augsburg) is a population-
30 based sample from the general population living in the region of Augsburg, Southern
31 Germany. The KORA F4 survey, a follow-up of the KORA S4 prospective cohort (1999-
32 2001), was conducted from 2006 to 2008 and included a total of 3080 participants. Clinical
33 and demographic information, as well as peripheral blood for 'omics' analyses, were
34 collected; details on the standardized examinations, interviews and tests conducted in the
35 KORA study have been previously described^{21, 22}. This study acted as discovery cohort in the
36 cross-sectional association study of plasma proteins and renal function (Fig. 1A).
37
38
39
40
41
42
43
44
45
46

47 Included in the replication phase were the Nord-Trøndelag Health Study (HUNT),
48 namely the third survey (HUNT3) from this population-based study from Norway with data on
49 participants of European descent²³; the INTERVAL Study (INTERVAL), a randomized trial
50 assessing blood donation practices across the UK with extensive phenotyping available for
51 50,000 participants of European descent²⁴; and the Qatar Metabolomics Study on Diabetes
52 (QMDIAB), a cross-sectional case-control study on type 2 diabetes from participants of Arab,
53
54
55
56
57
58
59
60

1
2
3 South Asian and Filipino descent in Qatar²⁵. Population characteristics from the four studies
4 are shown in Table 1. [Information on data availability is given in Supplemental Note 1.](#)

7 **Sample collection and proteomic profiling**

8
9 [EDTA plasma samples collected by the studies following standardized procedures were](#)
10 [centrifuged, aliquoted and stored at -80°C²⁶⁻²⁸.](#) [Samples for proteomic profiling and GFR](#)
11 [estimation were taken at the same time.](#)

12
13
14
15 [Proteomic profiling in all participating studies was done using SOMAscan](#)
16 [\(SomaLogic, Inc\), an aptamer-based, affinity proteomics platform^{10, 29-31}.](#) [Plasma samples](#)
17 [from KORA, HUNT3 and INTERVAL were shipped on dry ice to SomaLogic, Boulder, CO,](#)
18 [and proteomic profiling was performed using a SOMAscan panel of 1029 proteins for](#)
19 [KORA²⁶, 3,622 for INTERVAL²⁷ and 5000 for HUNT3²⁸.](#) [In the QMDiab cohort, the kit-based](#)
20 [SOMAscan platform was run by the Weill Cornell Medicine - Qatar \(WCM-Q\) proteomics core](#)
21 [following protocols and instrumentation provided by SomaLogic Inc., under supervision of](#)
22 [SomaLogic personnel, to measure 1,129 proteins in plasma samples²⁶.](#) [The samples were all](#)
23 [measured by individuals blinded to the identities corresponding to the samples.](#)

34 ***Measurement of plasma proteins***

35
36 ~~Proteomic profiling in all participating studies was done using SOMAscan (SomaLogic, Inc),~~
37 ~~an aptamer-based, affinity proteomics platform^{10, 29-31}.~~ In summary, fluorescently labeled
38
39 single-stranded synthetic nucleotides (Slow Off-rate Modified Aptamers, SOMAmers)
40
41 immobilized on streptavidin-coated beads are incubated with plasma samples to capture
42
43 proteins and generate SOMAmer-protein complexes. Washing steps eliminate unbound
44
45 SOMAmers and unbound/non-specifically bound proteins. The next steps are biotin-labeling
46
47 and photocleavage to liberate SOMAmer-protein complexes from the beads. This is followed
48
49 by incubation in a buffer disrupting nonspecific interactions, recapturing the biotin-labeled
50
51 protein/aptamer complexes in streptavidin-coated beads and additional washing steps to
52
53 remove nonspecific SOMAmers. These are then eluted from the target proteins and
54
55 quantified on custom DNA microarrays using deposited SOMAmer-complementary
56
57 oligonucleotides, which produces measurements in relative fluorescence units (RFU) as
58
59
60

proxies to protein concentrations. Quality control at the sample and SOMAmer levels using control aptamers and calibrator samples was conducted by the manufacturer. Based on standard samples included on each plate, the resulting raw intensities are processed using a data analysis work flow including hybridization normalization, median signal normalization and signal calibration to control for inter-plate differences.

In KORA F4, a random subset of 1,000 participants of those with omics data was selected for proteomic profiling with the SOMAscan assay featuring 1,129 protein-specific SOMAmer probes²⁶ in fasting plasma samples. SOMAscan QC resulted in the exclusion of 29 proteins and one sample, and five proteins were further excluded due to cross-reactivity (publicly available communication “SSM-064_Rev_0_DCN_16-263” issued by SomaLogic), producing data on $k = 1,095$ proteins in 999 participants in the discovery dataset (Supplementary Table 1). Proteomics data from INTERVAL featured 2,994 proteins in 3,301 participants²⁷, HUNT included 1,054 proteins in 2,432 individuals¹² and in-QMDiab 1,130 in 352 participants²⁶ after study-level quality control. for an overlap of 993 available proteins across datasets; further details on the proteomics profiling from the samples included in this study are described elsewhere^{12, 26, 27}. Protein mapping to several identifiers was provided by the manufacturer (Supplementary Table 1).

Outcome definitions

Our first analysis is a proteome wide association study: we investigated associations between proteins and renal traits as outcomes, using linear regression models with adjustment for potential confounders. The primary outcomes studied in this analysis were estimated glomerular filtration rate from serum creatinine (eGFR) and chronic kidney disease (CKD), given their availability in all included studies.

eGFR was calculated using the Chronic Kidney Disease Epidemiology Collaboration (CKD-EPI) equation with serum creatinine⁴ with the R package *nephro* v1.2³². Serum creatinine was measured using the modified kinetic Jaffé reaction in KORA, HUNT and QMDiab (and calibrated by multiplying by 0.95³³), and a nuclear magnetic resonance (NMR) platform (Nightingale Health) in the INTERVAL study. Pearson's correlation between serum

1
2
3 creatinine-based eGFR and this NMR-based eGFR variable was estimated in KORA (Suppl.
4 Fig. 1). Chronic Kidney Disease (CKD) was defined as eGFR <60 ml/min/1.73 m².³⁴.

5
6
7 We performed some analyses for outcomes available only in the discovery study.
8
9 Urinary albumin and urinary creatinine were used to calculate urinary albumin-to-creatinine
10 ratio (uACR) and its derived parameter microalbuminuria (MA, defined as uACR > 30 mg/g).
11 eGFR decline was defined as $\log(\text{eGFR})_{\text{follow-up}} - \log(\text{eGFR})_{\text{baseline}}$ divided by the follow-up
12 time, where KORA F4 (2006-2008) was used as baseline and KORA FF4 (2013-2014) its
13 follow-up survey. Sensitivity analyses were also run using eGFR_{cys} (derived from CKD-EPI
14 equation using cystatin C¹³).

21 **Definition of covariates**

22
23 Covariates used in the regression analyses were: age at the time of examination, sex, BMI,
24 smoking status, diabetes (yes/no), hypertension (yes/no), log-transformed triglycerides, high-
25 density lipoprotein, and intake of lipid-lowering drugs (yes/no). See Supplemental Note [4-2](#)
26 for precise cohort-specific definitions of covariates used.
27
28
29
30
31

32 **Statistical analysis**

33
34 Data preprocessing and statistical analyses were conducted using the R language for
35 statistical computing v.3.6.0³⁵. Prior to statistical analysis, proteomic data was log-
36 transformed and standardized. Linear regression was used to examine the association
37 between protein levels and continuous kidney traits (log-transformed eGFR, urinary-to-
38 albumin creatinine ratio, eGFR change), whereas logistic regression was used for binary
39 kidney traits (CKD, MA). Multiple testing was accounted for using a Bonferroni correction
40 considering the total number of investigated proteins at each stage (k = 1,095 in discovery).
41
42
43
44
45
46
47
48

49 Sensitivity analyses in the discovery sample included regression models with serum
50 creatinine-based eGFR as outcome and no adjustment for BMI or diabetes, as well as
51 models including cystatin C-based eGFR as outcome and the same set of covariates from
52 the main model. Pearson's correlations between the regression coefficients resulting from
53 the sensitivity and the main analyses were calculated. [Interaction analyses were also](#)
54
55
56
57
58
59
60

1
2
3 conducted for the proteins identified at discovery by adding an interaction term, considering
4 age, sex and smoking as interactors, to the fully adjusted model (Supplemental Note 3).
5
6

7 For those protein-outcome pairs significantly associated in the discovery, two
8 replications were conducted: a European replication (R1) and a replication in an admixed
9 population (R2) (Fig. 1A). Replication was defined as $p < 0.05$ and consistent direction of
10 effect as in the discovery study. The European replication for eGFR consisted of the meta-
11 analysis of results from HUNT and INTERVAL using ~~a modified form of~~ the Stouffer's method
12 ~~(weighted by the sample size of each individual study)³⁶, as common meta-analysis~~
13 ~~approaches using regression coefficients could not be used due to the relative nature of the~~
14 ~~SOMAscan data (measured in relative fluorescence units, RFU).~~
15 a P-value combination method especially useful when raw data cannot be pooled across
16 studies – which is the case with aptamer-based measurements, where data in relative
17 fluorescence units (RFU) is not directly comparable across studies. Stouffer's method, also
18 known as “inverse normal” or weighted Z-test, is a P-value combination method taking the P-
19 values for the i-th study (p_i), transforming them by the inverse normal distribution function
20 and weighing them using the square root of the sample sizes as weights (w_i). The sum is
21 then computed, and the combined P-value is obtained using the distribution of the resulting
22 statistic, $T = \sum w_i H(p_i)$ ³⁶.
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40

41 For CKD, only the HUNT study was used in the European replication (INTERVAL had
42 only one case), and the admixed population replication was based on the results of QMDiab.
43 Our final set of trans-ethnic associations (R3) were those pairs of proteins-outcomes that
44 were replicated in both R1 and R2. Replicated eGFR-associated proteins were taken to the
45 next stages of the analysis: proteomic target validation, enrichment analyses and Mendelian
46 randomization.
47
48
49
50
51
52

53 **Validation of proteomic targets**

54 We examined the plasma levels of proteins measured by Proximity Extension Assay (PEA)
55 technology (Olink) in a subgroup of randomly selected participants from the KORA F4 study
56 (N = 173) (³⁷, DOI: 10.1021/acs.jproteome.0c00641). In brief, protein abundance was
57
58
59
60

1
2
3 quantified by real-time PCR using the PEA proteomic technology (Olink), producing relative
4 quantification data reported in NPX values (normalized protein expression levels, on log₂
5 scale); NPX values were intensity normalized with the plate median for each assay as the
6 normalization factor, and samples and proteins that did not pass the quality control were
7 excluded (DOI: 10.1021/acs.jproteome.0c00641). Seven of the most relevant proteins (i.e.
8 cystatin-C, RELT, IGFBP-6, myoglobin, TNF sR-I, RGMB and FSTL3), the novel proteins
9 here reported (i.e. JAM-B and contactin-4), as well as three of the proteins identified in the
10 causal inference analysis (i.e. carbonic anhydrase 3, MIA, and cystatin M) were included in
11 this subset of proteomic measurements. Of note, testican-2 was not measured in this assay.
12 Scatterplots of the aptamer-based and PEA measurements, annotated with their Pearson's
13 correlation and statistical significance, are shown in Suppl. Fig. 2.

14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

Information on specificity and cross-reactivity of the aptamers was available from three independent studies^{27, 38, 39} for 54 of the 57 proteins identified to be trans-ethnically associated with eGFR. Target specificity issues (i.e. comparable binding observed to a target that is not the product of the same gene) was-were observed in four cases (ephrin-A5, IGFBP-5, hemojuvelin, and cystatin SA)^{27, 38} (Supplementary Tables 2-4). Moreover, in previous studies, 23 of the 57 proteins were directly validated via mass spectrometry in blood plasma/serum, and other biological matrices³⁹ (Supplementary Table 2), and 49 using solution affinity measurements^{27, 38} (Supplementary Tables 3-4).

Functional annotation, and enrichment and expression analysis analyses

Annotation was done using the R package *InterMineR* v1.6.1⁴⁰, a tool facilitating access to data from the HumanMine release 6.0 (May 2019). DAVID v.6.8⁴¹ was used to look for annotations for Gene Ontology Terms (molecular function, biological process), and pathway and gene information, as well as to identify publications relevant to the set of 57 replicated proteins. Gene information was retrieved from the human assembly GRCh37 (hg19) using BioMart v.4⁴².

1
2
3 To investigate the expression patterns of the 57 eGFR-associated proteins and their
4 corresponding protein coding genes across tissues, we used proteomics and RNA-seq
5 expression data from the ProteomicsDB^{43, 44} and the Genotype-Tissue Expression (GTEx)
6 database⁴⁵. The data presented and described in this manuscript were generated on Oct. 2,
7 2020 through a multi-gene query on the ProteomicsDB Analytics Toolbox portal from:
8 <https://www.proteomicsdb.org/proteomicsdb/#analytics/expressionHeatmap> and GTEx portal
9 <https://www.gtexportal.org/home/multiGeneQueryPage>

Protein-protein interaction (PPI) network analysis

10
11 We queried STRING⁴⁶, the protein-protein interaction server, to examine the relationship
12 between the proteins that were identified as robustly associated with eGFR across studies
13 and ethnicities (k = 57 trans-ethnically eGFR-associated proteins). We used the set of
14 SOMAscan proteins available across studies as background (k = 993), adding no additional
15 interactors (proteins) to the network during the analyses, and considered a minimum required
16 interaction score for a medium confidence (0.400).

Mendelian randomization

17
18 Mendelian Randomization (MR), an instrumental variable method used to infer causality,
19 leverages the natural randomization inherent in the (random) assortment of genes during
20 gamete formation to assess the effect of lifelong exposures on health outcomes⁴⁷. Single-
21 nucleotide polymorphisms (SNPs) are used as instrumental variables (or instruments), given
22 that their alleles are randomly assigned to individuals prior to any exposures/outcome and
23 that they are non-modifiable, thus minimizing the risk of reverse causation and confounding
24⁴⁷. The idea behind MR is that if genetic variation produces differences mirroring the
25 biological effects of environmental exposures that alter disease risk, then genetic variation
26 itself should be related to disease risk by having an influence on the exposure^{47, 48}. MR uses
27 SNPs as surrogates for an exposure of interest, allowing the estimation of the effects of life-
28 long, genetically determined “exposures” on health outcomes⁴⁷. MR produces robust causal
29 inference estimates if the SNPs used are valid instruments – that is, if they meet the three
30 assumptions upon which MR relies: SNPs must be strongly associated with the exposure,
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

1
2
3 and not associated with either (measured or unmeasured) confounders or with the outcome
4 (only potentially through the exposure) ⁴⁹. Causality in MR is thus defined as the modification
5 of an exposure leading to a change in the outcome, where the inferred causal effects by MR
6 do not necessarily imply the existence of a straightforward interpretation with respect to
7 direct causal factors ⁴⁸.

8
9
10
11
12
13 To investigate whether genetic liability to lower or higher eGFR causally alters plasma
14 protein levels and vice versa, MR was conducted ~~Causality was assessed~~ in the set of 57
15
16 proteins whose associations with eGFR showed trans-ethnic replication. Two-sample
17
18 bidirectional Mendelian randomization ¹⁹ (MR) was used to ~~estimate~~ infer the causal effect of
19
20 renal function (eGFR as proxy thereof) on plasma protein levels (forward MR) and vice versa
21
22 (reverse MR, Fig. 1B). Results from publicly available genome-wide association studies
23
24 (GWAS) for (a) eGFR from the CKDGen consortium (meta-analysis of European-ancestry
25
26 populations ⁵⁰), and (b) plasma proteins from INTERVAL²⁷ and AGES-Reykjavik³⁹ were used
27
28 to perform MR using MRBase⁵¹. A detailed account on the MR methods, data sources and
29
30 analyses conducted is available in Supplemental Note 34.

31 32 33 34 ***Instrument selection***

35
36 In the forward MR (i.e. assessing the effect of renal filtration on protein levels), 256 SNPs
37
38 associated with eGFR at genome-wide significance in the CKDGen results were selected as
39
40 candidate instrumental variables (IV). These SNPs were then filtered based on their
41
42 relevance to renal function (associated with BUN, a complementary renal trait, with an
43
44 opposite direction of effect, $N_{IV} = 47$) and clumped based on linkage disequilibrium ($r^2 = 0.01$
45
46 and $Kb = 10,000$) to identify independent variants ($N_{IV}=41$). Summary statistics on 41 SNP-
47
48 eGFR associations were extracted from the CKDGen results, and its corresponding SNP-
49
50 protein associations were extracted from the INTERVAL results for 47 proteins. For
51
52 investigating the causal effects of eGFR on proteins, 47 eGFR-protein relationships were
53
54 instrumented by 41 SNPs.

55
56
57 For the reverse MR (i.e. interrogating the causal effect of proteins on renal filtration),
58
59 genome-wide significant cis-SNPs for 22-28 proteins were identified in the INTERVAL results
60

1
2
3 as candidate IV and LD clumped (same criteria as forward MR). Summary statistics on SNP-
4 protein associations for 24-28 proteins were extracted from the INTERVAL results, and its
5 corresponding SNP-eGFR associations were extracted from the CKDGen results. The same
6 strategy was followed to identify instruments in the AGES-Reykjavik results; ~~for the 36~~
7 ~~proteins not found in INTERVAL~~; SNP-protein results were extracted from this dataset for
8 ~~seven 29~~ proteins, and ~~their corresponding~~ SNP-eGFR results were extracted for 26 proteins
9 from the CKDGen data. Further details of the genetic instrument selection and data
10 harmonization process are shown in Suppl. Fig. 23 and Supplementary Table 5. Thus for
11 investigating the causal effect of proteins on eGFR, 28-35 protein-eGFR relationships were
12 instrumented by 1-3-5 SNPs, of which 17 proteins were examined using data from both
13 INTERVAL and AGES-Reykjavik (Suppl. Fig. 4). ~~Further details of the genetic instrument~~
14 ~~selection and data harmonization process are shown in Suppl. Fig. 2 and Supplementary~~
15 ~~Table 5.~~

Data harmonization, phenotypic variance explained and instrument specificity

16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
Details on data harmonization, the handling of palindromic SNPs and calculating the
phenotypic variance explained by the SNPs are given in Supplemental Note 34. Harmonized
datasets used in the MR analyses are available in Supplementary Table 6.

In order to look for further evidence of horizontal pleiotropy, association between our
SNPs and other traits were searched for in the GWAS Catalog⁵² (Suppl. Table 7).

MR and sensitivity analyses

The primary MR analysis used inverse variance weighted (IVW) regression. In this method
the coefficient of the gene-outcome association is regressed on the coefficient of the gene-
exposure association with the intercept constrained to zero, assuming no directional
pleiotropy^{53, 54}. Since IVW requires two or more SNPs, in cases where only one SNP
instrumented the analysis, Wald's ratio (coefficient of the gene-outcome association divided
by the gene-exposure association) was calculated instead⁵⁴.

For MR analyses instrumented by more than two SNPs, three further MR methods were
used as sensitivity analyses⁵⁵. MR-Egger regression was used to assess pleiotropy, as this

1
2
3 method allows for horizontal pleiotropy and provides an estimate of the unbalanced
4 horizontal pleiotropic effects in its intercept⁵⁶. Weighted median⁵⁷ and weighted mode MR⁵⁸,
5 methods less sensitive to the presence of invalid instruments and to pleiotropic SNPs
6 behaving as outliers, were also used. A number of additional analyses were run to check for
7 outliers, directional pleiotropy and heterogeneity, as recommended^{55, 59}. Details are given in
8 Supplemental Note [34](#).

9
10
11 Causal estimates were assessed at a Bonferroni-corrected significance level, i.e. 0.05
12 divided by the number of proteins assessed in each MR direction (47 in forward and [28-51](#) in
13 reverse MR). Causal effects were considered robust if they were significant at Bonferroni $p <$
14 0.05 in the IVW or Wald estimator, and results from the [pleiotropy-robust](#) sensitivity MR
15 analyses examined to test for violations to MR assumptions.

16 **Expression analyses in [human kidney tissue cohort of CKD patients](#)**

17
18 The correlation between gene expression analysis of *SPOCK2*, [one of the genes](#) coding for
19 [the proteins](#) with ~~the strongest and most robust~~ evidence from the causal analysis (~~testican-~~
20 ~~2~~), and estimated glomerular filtration rate (eGFR), was [initially calculated for with](#)
21 ~~independent~~ data from microdissected tubulointerstitial components of human renal biopsies
22 from 26 individuals with CKD at different disease stages (I-IV)⁶⁰ (GEO accession:
23 GSE69438). [Gene expression of the protein-coding genes identified in MR \(*SPOCK2*, *CA3*,
24 *CST6*, *MIA*\) and renal traits was further assessed in \(a\) data from *Nephroseq* v5 \(N = 458\), a
25 \[platform of comprehensive renal disease gene expression datasets\]\(#\)⁶¹, and \(b\) human kidney
26 \[tissue resource characterised by RNA-sequencing \\(N = 427, see Supplemental Note 5\\)\]\(#\)¹.](#)

27
28 [Within *Nephroseq*, univariate correlation analyses between eGFR and gene](#)
29 [expression were conducted separately in study pre-defined histological compartments of the](#)
30 [human kidney \(i.e. glomerular and tubulo-interstitial\) in 458 available kidney samples from](#)
31 [three datasets of patients with kidney disease \(Ju et al.⁶⁰, Sampson et al.⁶², and Reich et al.](#)
32 [63\), and one dataset of “apparently” healthy renal tissue \(Rodwell et al.⁶⁴\). The correlations](#)
33 [were meta-analysed using inverse variance weighted meta-analysis with random effects](#)
34 [models](#)⁶⁵, [and heterogeneity was assessed using Cochran’s Q test.](#)

Multivariable regression analyses were conducted in the human kidney resource (N = 427) described in ¹. In brief, we constructed linear regression models with renal expression of each candidate as the response variable; while eGFR and histologically confirmed measures of structural kidney damage were used as dependent variables together with age, sex, body mass index, 3 genetic principal components, diabetes and a variable number of surrogate variables (29 for eGFR and 26 for all histology phenotypes) ^{1, 66}. eGFR estimation was based on circulating levels of creatinine, as reported before ¹. Histologic measures of structural integrity (glomerular sclerosis, glomerular Bowman's capsule thickening, tubular atrophy, interstitial fibrosis, interstitial inflammation and vascular lesions) were assessed microscopically and scored on a semi-quantitative scale (whereby 0 indicates no or minimal damage and 3 is consistent with the highest degree of structural injury), as reported before ⁶⁷.

Results

Figure 1 illustrates the design of the present study. First, a cross-sectional association study was performed to identify proteins associated with renal function parameters in a discovery-replication setting: KORA F4 acted as the discovery, and INTERVAL, HUNT3 and QMDiab as replication studies (Fig. 1A). Replicated trans-ethnic protein associations were then assessed for causality using two-sample Mendelian randomization, using data from the largest genome-wide association (GWA) studies available for the traits of interest (CKDGen, INTERVAL and AGES-Reykjavik) (Fig. 1B).

Cross-sectional association of plasma proteins and renal function

Population characteristics of the four cohorts included in the cross-sectional association study are shown in Table 1. The largest differences between the discovery and the replication studies were observed in the age distributions, smoking habits, blood lipid levels, eGFR distribution, and CKD/diabetes prevalence.

Results from discovery study

The association between 1,095 plasma proteins and eGFR/CKD was assessed in the KORA F4 study (N = 995). A total of 80 proteins were significantly associated with eGFR ($p < 0.05/1095$). The top 3 negative associations (i.e. higher eGFR associated with lower plasma

1
2
3 protein levels) were observed with cystatin C ($\beta = -0.068$ [95% CI = -0.078,-0.059] change in
4 log-transformed eGR per standard deviation increase in protein level, $p = 2.63E-40$), tumor
5 necrosis factor receptor superfamily member 19L (REL1; $\beta = -0.063$ [-0.073,-0.053], $p =$
6 $7.82E-33$) and beta-2-microglobulin (b2-microglobulin; $\beta = -0.059$ [-0.070,-0.050], $p = 6.16E-$
7 30), whereas the strongest positive association (i.e. higher eGFR associated with higher
8 plasma protein levels) was that of testican-2 ($\beta = 0.036$ [0.026, 0.045], $p = 2.066E-13$)
9 (Supplementary Table 1, Suppl. Fig. [3-5](#) and [4A6A](#)). Of note, 34 of these 80 proteins were
10 also associated with CKD (Supplementary Table 1).

11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
Sensitivity analyses showed that 71 of the 80 eGFR-associated-proteins identified in
the main analysis were consistently associated with cystatin C-based eGFR, with a high
correlation between regression coefficients ($r = 0.841$, $p < 2.2E-16$). Models with no
adjustment for BMI or diabetes produced highly similar estimates to those obtained in the
main analysis ($r = 0.99$, $p < 2.2E-16$ for both; Supplementary Table 8). Likewise, the
exclusion of individuals with CKD (N = 38) did not significantly affect the correlation between
the plasma levels of proteins and log-transformed eGFR (Suppl. Fig. 7). Interaction analyses
suggested the negative associations between log-transformed eGFR and five plasma
proteins (b2-Microglobulin, IGFBP-6, FSTL3, JAM-B, and renin) were accentuated with age
(i.e. each additional year of age made the association stronger) (Supplementary Table 9).

41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
To further explore proteomic associations with renal function, additional renal
outcomes were assessed in the discovery cohort. eGFR change was associated with five
proteins, whereas three proteins were negatively associated with uACR (epidermal growth
factor receptor (EGFR), mitochondrial superoxide dismutase [Mn] (SOD2) and coagulation
factor X (F10)). No proteins were significantly associated with microalbuminuria
(Supplementary Table [910](#)).

Results from replication studies

Serum creatinine was the only available trait across all replication cohorts (Fig. 1A), thus only
associations with eGFR/CKD were further explored.

1
2
3 The European replication (R1) was done using HUNT3 and INTERVAL (Suppl. Fig.
4 [4B6B-C](#)); results from this analysis confirmed the association of 62 of the 76 proteins
5 available across studies. The second replication round (R2) was performed in QMDiab
6 (Suppl. Fig. [4D6D](#)), a population of admixed-ancestry; this confirmed 63 eGFR-protein
7 associations. [High correlation between z-values from the discovery and replication studies](#)
8 [was observed \(correlations ranging from 0.66 with INTERVAL to 0.93 with HUNT3, Suppl.](#)
9 [Fig. 8\)](#). The overlap of the proteins replicated in R1 and R2 was the final set of 57 robustly
10 replicated trans-ethnic eGFR-protein associations (R3, Supplementary Table [4011](#)). Figure 2
11 shows the cross-sectional effect estimates for the top 10 protein-eGFR associations across
12 the four cohorts; INTERVAL, a largely healthy and younger population, showed the smallest
13 effect sizes, whereas the strongest effects were observed in HUNT3, a cohort of older
14 individuals with lower mean eGFR and higher CKD prevalence. Two novel proteins were
15 identified (contactin-4 and junctional adhesion molecule B, JAM-B), and all of the 57 proteins
16 were replicated in the cystatin C-based eGFR sensitivity analysis (Supplementary Table 8).

17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32 All 34 CKD-protein associations from the discovery phase were replicated in HUNT
33 and 20 replicated in QMDiab; these 20 were thus considered trans-ethnically robust
34 (Supplementary Table [4112](#)). Figure 3 shows the sets of replicated proteins associated with
35 eGFR/CKD. JAM-B, one of the novel proteins described here, is one of the 19 proteins
36 associated with both eGFR and CKD.

37 ***Functional annotation, and enrichment and expression analysis analyses***

38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
An extended annotation file including information on GO terms, pathways (KEGG,
BIOCARTA), and protein domain annotation (INTERPRO) is given in Supplementary Table
[4213](#), and the retrieved gene information in Supplementary Table [4314](#). Although several
pathways, biological processes and molecular functions were represented, no enrichment for
Gene Ontology (GO) terms or pathways was observed.

[Protein abundance based on quantitative mass spectrometry assays \(from](#)
[ProteomicsDB\) across different tissues, cell lines, and body fluids is shown in Suppl. Fig. 9;](#)

1
2
3 peptides for 15 of the 57 eGFR-associated proteins were detected in multiple human tissues
4
5
6 and body fluids (including kidney tissue); of these ubiquitously expressed proteins, beta-2
7
8
9 microglobulin and glutathione S-transferase Pi are most notable (Suppl. Fig. 9). Most protein-
10
11
12 coding genes demonstrated ubiquitous expression across virtually all human tissues,
13
14
15 including kidney tissue, represented in the ProteomicsDB datasets, where genes like *B2M*
16
17
18 and *GSTP1* were highly expressed across tissues in both RNA-seq and microarray
19
20
21 expression datasets (Suppl. Fig. 10). GTEx tissue expression data was largely concordant
22
23
24 with the observations from ProteomicsDB (Suppl. Fig. 11).

Protein-protein interaction (PPI) network analysis

25
26
27
28
29
30 We queried STRING⁴⁶, the protein-protein interaction server, to examine the relationship
31
32 between the 57 replicated proteins associated with eGFR. The obtained network features 56
33
34 nodes representing proteins, connected by 60 edges representing predicted functional
35
36 protein-protein associations, with an average number of interactions of 2.14 at a minimum
37
38 interaction score of 0.400 (medium confidence) in the network (Suppl. Fig. 12). The expected
39
40 number of edges was 36 (PPI enrichment p-value: 1.12E-04). This means the examined set
41
42 of 57 eGFR-associated proteins have more interactions among themselves than what would
43
44 be expected from a random set of proteins drawn from the set of proteins included in the
45
46 SOMAscan studies, and thus suggest these proteins might be meaningfully biologically
47
48 connected.

Mendelian randomization

49
50
51
52
53
54 To assess whether genetic susceptibility to higher or lower plasma levels of the 57 proteins
55
56 identified as trans-ethnically associated with eGFR may affect this renal trait, and whether
57
58 genetic liability to altered eGFR causally alters circulating levels of plasma proteins, might be
59
60

causally affected by or involved in changes in renal function, two-sample bidirectional Mendelian randomization (MR) was conducted (Fig. 1B).

Forward MR: eGFR has an effect on testican-2

In the forward direction of the MR (i.e. ~~assessing~~ inferring the effect of eGFR on levels of 47 proteins), 40 SNPs explaining 1.59% of the variance of eGFR (Supplementary Table ~~14~~15) were used as genetic instruments.

Plasma levels of Seven-7 proteins were identified as causally affected by eGFR according to the IVW model (cystatin M, cathepsin H, ephrin type-B receptor 6, insulin-like growth factor-binding protein 6, testican-2, melanoma-derived growth regulatory protein, and netrin receptor UNC5C; Supplementary Table ~~15~~16 and Suppl. Fig. 13-19). This means that if eGFR is somehow altered by means of an intervention mimicking the effect of the SNP on eGFR, MR suggests the plasma levels of these proteins will change in the predicted direction. Although no evidence of directional pleiotropy, particularly influential SNPs or instrument heterogeneity was observed (with the exception of IGFB6, Suppl. Tables ~~16~~17-18), the pleiotropy-robust sensitivity MR analyses did not provide further evidence of support causality for six of them (i.e. non-significant results), raising the possibility of suggesting the IVW findings might be being driven by undetected balanced-horizontal pleiotropy⁶⁸. In contrast, the causal effect of eGFR on testican-2 was supported by more than one MR method (weighted median $\beta = 6.104$ [95% CI, 2.808 to 9.401] unit increase in testican-2 levels per unit increase in log-transformed eGFR, $p = 2.84E-04$) (Table 2 and Figure 4). Furthermore, there was no evidence of outliers, influential SNPs or instrument heterogeneity (Suppl. Fig. 517).

Eleven of the SNPs instrumenting the forward MR analysis were identified as potentially pleiotropic (associations at $p < 5E-08$ with other related traits, Supplementary Table 7). Results from the restrictive MR conducted after excluding these pleiotropic variants were not statistically significant, but in agreement with those from the main analysis in terms of their direction and size of effect (Suppl. Table 1920).

Reverse MR: MIA, cystatin M and carbonic anhydrase III affect eGFR

1
2
3 In the reverse direction of the MR (i.e. assessing the effect of 28-35 proteins on eGFR),
4 one to three-five cis-SNPs explaining 0.914-02%-29.33% of phenotypic variance were used
5 as genetic instruments (Supplementary Table 1415). Of note, no causal effect of testican-2
6 on eGFR was identified in this direction of the MR (Wald's ratio $p = 0.053$ for variant
7 $rs1245547$ from AGES-Reykjavik and $p = 0.06$ for variant $rs1245540$ from INTERVAL)
8 (Suppl. Table 16).

9
10
11 A negative effect of melanoma inhibitory activity (MIA) on eGFR (i.e. higher protein levels
12 have a negative effect on eGFR) was identified in the fixed-effect IVW and weighted-median
13 models using the 3 SNPs identified in the INTERVAL pGWAS (Table 2, Suppl. Fig. 6A20A-
14 B), meaning MR suggests that if plasma protein levels are lowered by means of an
15 intervention mimicking the effect of the SNP on MIA, eGFR will increase. No evidence of
16 influential SNPs was observed, yet the funnel plot suggested directional pleiotropy
17 (Supplementary Tables 1617-1819, Suppl. Fig. 206C). One SNP instrumenting this analysis
18 was identified as potentially pleiotropic (Supplementary Table 7). The causal effect estimated
19 using AGES-Reykjavik pGWAS data was in agreement with this MR estimate (Suppl. Table
20 16).

21
22
23 A pPositive effects of carbonic anhydrase III and cystatin M on eGFR was-were also
24 identified (Wald's ratio $p = 5.04E-04$ and $8.41E-05$, respectively, with INTERVAL pGWAS data)
25 (Table 2). Although further sensitivity analyses were not conducted given there was only one
26 cis-SNP available for each protein, no gene-trait associations were found for these SNPs in
27 the GWAS Catalog⁵², suggesting pleiotropic effects to be unlikely. The effect of carbonic
28 anhydrase III on eGFR estimated using pGWAS data from AGES-Reykjavik was in line with
29 the aforementioned estimate (Suppl. Table 16).

SPOCK2 expression in renal samples from CKD patients

Gene expression in kidney tissue

~~A moderate non-significant~~ An initial assessment of the expression of the gene coding for *testican-2*, *SPOCK2*, in tubulointerstitial components of human renal biopsies from 26 individuals with CKD⁶⁰ found no correlation was found between expression of *SPOCK2* and with eGFR at the time of biopsy (Pearson's $r^2 = 0.25$, $p = 0.22$) ~~on tubulointerstitial components of human renal biopsies from 26 individuals with CKD⁶⁰~~ (Suppl. Fig. 721). Further univariate analyses conducted with the *Nephroseq* data showed a statistically significant correlation between eGFR and *SPOCK2* gene expression in glomerular compartment/kidney cortex ($r = 0.242$, $p = 0.033$) (Suppl. Table 21). We then conducted additional analyses using RNA-sequencing-characterised human kidney transcriptomes from a resource with up to 427 individuals with matching gene expression and renal phenotype data (as described before in ^{1, 66, 67} and Suppl. Table 22). There was no significant association between eGFR and renal *SPOCK2* expression in this dataset (N=427, Supplemental Note 5; Suppl. Table 23). However, in the subset with information on histology phenotypes (N=283), *SPOCK2* expression was negatively associated with both tubular atrophy ($\beta = -0.094$ [-0.177, -0.011], $p = 0.03$) and interstitial fibrosis ($\beta = -0.093$ [-0.175, -0.010], $p = 0.03$), and *CST6* expression was negatively associated with glomerular sclerosis ($\beta = -0.094$ [95% CI = -0.174, -0.013], $p = 0.02$).

Discussion

We conducted an association study of plasma proteomics and eGFR/CKD following a discovery-replication approach involving four independent studies. We confirm known protein associations and identify two novel potential biomarkers. Furthermore, we performed two-sample bidirectional Mendelian randomization to ~~identify infer causal relationships causality in between~~ the eGFR-associated protein associations, and found evidence of causality underlying four eGFR-protein associations.

Eighty proteins were found to be associated with eGFR in our discovery analysis, with trans-ethnic replication confirming 57 of these. Nineteen of these proteins were also found to be associated with CKD. Although our analyses use serum-creatinine-based eGFR due to its

1
2
3 availability across studies and its utility in clinical practice, models using cystatin C-based
4 eGFR show the results are robust to GFR estimation method. Likewise, further sensitivity
5 analyses indicate the associations are largely independent of adjustment for BMI or diabetes.
6
7

8
9 Interestingly, from the five proteins with a significant interaction with age, protein trajectories
10 changing with age have been reported for FSTL3, IGFBP6, JAM-B and renin ⁶⁹.

11
12
13 We identify several well-known biomarkers of renal function^{10, 16, 17, 70, 71}, supporting
14 their current use as kidney function biomarkers and the validity of our analyses. Our results
15 are also in line with those reported by other renal studies using the same aptamer-based
16 platform and similar proteomic profiling technologies^{10, 16-18}: we replicate 15 of the proteins
17 identified in the pioneer SOMAScan study of plasma from 42 CKD patients¹⁰, five of the
18 proteins associated with lower baseline eGFR and 5-year eGFR decline in a study examining
19 80 circulating proteins in ~1,000 participants¹⁶, and a large number of proteins reported in a
20 recent SOMAScan study of 2,893 plasma proteins in 389 Swedish individuals¹⁷. Moreover,
21 five of our proteins (TNF SR-I and -II, TAJRELT, CD55 and CCL14) were included in a
22 signature capturing the inflammatory process underlying end-stage renal disease in diabetic
23 cohorts⁷¹, and another five proteins (b2-microglobulin, cystatin C, DAF, MP2K2 and testican-
24 2) are found in a set of were included proteins in a meant to reflect renal health in a “stand-
25 alone” test meant to reflect renal health ³⁸. Interestingly, 40% of our proteins were identified
26 in podocyte exosome-enriched urine, suggesting their involvement in cellular functional
27 processes underlying glomerular filter permeability⁷². A recent aptamer-based study found
28 126 proteins associated with baseline eGFR ¹⁸ – an overlap of 43 proteins independently
29 reported in both their and our study, including well-known proteomic biomarkers (e.g.
30 cystatin-C, b2-microglobulin) and testican-2, is listed in Suppl. Table 11. The protein and
31 gene expression results for our 57 replicated proteins from ProteomicsDB and GTEx
32 databases confirm the renal expression of these proteins. Furthermore, the ubiquitous
33 expression across tissues and evidence from the PPI network analysis suggest their
34 involvement in common functional pathways.
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

1
2
3 Among the 57 eGFR-associated proteins we identified, we identify two have never
4 been reported in association with eGFR nor any other renal trait novel eGFR-associated
5 proteins: contactin-4, a protein involved in neuronal network development and plasticity; and
6
7 junctional adhesion molecule B (JAM-B), involved in cellular junctions in epithelial tissue. Of
8
9 note, the latter was also recently identified in the study by Ngo and colleagues ¹⁸, with
10
11 consistent direction of effects to our findings. Glomerular filtration of proteins (determined by
12
13 their molecular weight, molecular charge and shape, and their interactions with other
14
15 molecules) influences their plasma levels^{73, 74}. Moreover, proteins detected in urine and
16
17 plasma are derived mainly from glomerular filtration of plasma proteins and epithelial cell
18
19 secretion of soluble proteins^{73, 74}. Contactin-4 (113.45 kDa) and JAM-B (33.21 kDa) are
20
21 present in plasma but not in urine, suggesting they are either not filtered at the glomerular
22
23 capillaries (perhaps due to interactions with other proteins) or filtered but later reabsorbed
24
25 into blood from the tubules; this in turn suggests changes in their plasma levels may reflect
26
27 changes in glomerular and tubule function⁷³. Moreover, the direction of effect observed for
28
29 JAM-B (higher levels significantly associated with both lower eGFR and CKD OR > 1) is
30
31 supportive of this protein being related to decreased glomerular filtering.
32
33

34
35 We also identified proteins associated with complementary renal phenotypes (eGFR
36
37 change, uACR and MA). eGFR decline was associated with five proteins, of which neurexin-
38
39 1-beta is the only not having been previously described as related to eGFR ^{10, 16, 17}. Of these
40
41 five proteins, DAN, TNF sR-1 and FSTL 3 were also identified as negatively associated with
42
43 eGFR decline in ¹⁸. There was no overlap between the set of proteins associated with uACR
44
45 in our study and the proteins identified in a similar study⁷⁵, which may be explained by the
46
47 albuminuria samples and eGFR being measured at different time points in their study⁷⁵. No
48
49 proteins were significantly associated with MA, possibly due to its low prevalence (5.9%,
50
51 Table 1) in KORA.
52
53

54
55 To investigate whether genetic susceptibility to renal function (using eGFR as a proxy
56
57 thereof) or plasma protein levels may have a causal effect on the other, and identify potential
58
59 causal disease pathways represented by proteins, better understand the biological
60

1
2
3 ~~mechanisms in which these proteins might be involved~~, bidirectional MR was performed to
4 identify causal effects of proteins on eGFR and vice versa. To date, only one study using MR
5 to assess causality of kidney function and proteomic biomarkers has been reported⁷⁶;
6
7 however, its focus, studied population and biomarker selection are markedly different to ours.
8
9

10
11 Our forward MR analyses initially suggested ~~the existence of a causal effect of~~ renal
12 function ~~may affect on the protein~~ plasma levels of seven proteins, although sensitivity
13 analyses seem to indicate horizontal pleiotropy may be at play for six of them, given that their
14 estimates were compatible with the null (i.e. no causal effect). A ~~robust~~ causal ~~effect~~
15 ~~association of between~~ renal function ~~on and plasma levels of~~ testican-2 was
16 ~~observed/inferred, given with the~~ supporting evidence ~~offered by from pleiotropy-robust~~
17 sensitivity MR methods allowing for violations ~~to of~~ MR assumptions. Although consistent
18 with the main analysis, the results from the restrictive MR were not significant perhaps due to
19 reduced statistical power deriving from (a) fewer SNPs instrumenting the analysis, and (b)
20 the exclusion of SNPs that might be on the actual causal pathway of interest.
21
22
23
24
25
26
27
28
29
30
31

32 Testican-2 is a secreted protein of the SPARC family⁷⁷, a group of matricellular proteins
33 (MCPs) regulating extracellular matrix (ECM)-cell interactions and ECM processing⁷⁸, and is
34 involved in a number of biological processes (Supplementary Table [2024](#)). Given its
35 glomerular filtration and detection in urine⁷⁹, changes in its plasma levels may reflect
36 ~~changes in kidney function glomerular filtration alterations~~⁷⁴, ~~given its renal release into the~~
37 ~~bloodstream~~¹⁸. Interestingly, higher testican-2 plasma levels have also been associated with
38 ~~less eGFR loss over time and reduced odds of incident CKD~~¹⁸. Its protein-coding gene,
39 *SPOCK2*, is associated with both normal maintenance of organ and tissue integrity
40 (glomerular remodeling), as well as with wound healing and other responses to injury⁸⁰.
41 Enriched in human glomeruli in comparison to tubuli samples^{81, 82} and other non-renal
42 tissues⁸³, *SPOCK2* has been reported as a glomerular and podocyte-specific gene^{82, 84}.
43 ~~where recent evidence from immunohistochemistry and immunofluorescence of human~~
44 ~~kidney tissue confirm its glomerular expression, and podocyte-specific expression in adult~~
45 ~~human kidney samples at single-cell resolution~~¹⁸. ~~The statistically insignificant correlation~~
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

1
2
3 between eGFR and *SPOCK2* expression in tubules in our study might be thus explained by
4 the glomerular specificity of this gene. Moreover, *SPOCK2* has been reported as one of the
5 top downregulated genes in glomeruli from patients with diabetic kidney disease⁸¹, which
6 further suggests its involvement in glomerular renal disease. The association between eGFR
7 and *SPOCK2* renal expression did not reach the level of nominal statistical significance in
8 some experiments and this may be explained at least partly by low statistical power (i.e. in
9 the case of the dataset of 26 CKD patients) but the directionality of the association was
10 consistent across datasets. Moreover, the multivariable analyses showed higher scores of
11 histologic measures of renal structural damage to be negatively associated with *SPOCK2*
12 renal expression, and positively associated with *CST6* renal expression - the direction of
13 these associations was again consistent with that of the effects in the cross-sectional and
14 MR analyses. Indeed All in all, the agreement between the positive association cross-
15 sectional results reported by us here and by others¹⁸, as well as our and-MR results findings
16 and the association with histologic measures, indicate testican-2 and its protein-coding gene
17 *SPOCK2* may have an active role in renal health support the notion. Low of low plasma
18 levels of testican-2 may thus be as indicative of poor renal function, suggesting testican-2
19 this protein to be is a physiological biomarker of kidney health and disease progression,
20 rather than only as a filtration marker^{17, 18}. Although the reverse direction of this causal
21 association could not be tested in our analyses was not significant in our MR analyses, a
22 causal effect of testican-2 on renal function is also possible biologically plausible, given the
23 role MCPs might play in the shift from constructive ECM repair to tissue stiffening and
24 fibrosis^{78, 85}, and the recent evidence of its potential *in vitro* effects on human glomerular
25 endothelial cells (HGECs) motility¹⁸. Nevertheless, whether the utility of testican-2 as a
26 biomarker is related to its potential functional effects, and mechanisms influencing its blood
27 levels, requires further study¹⁸.

28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
Three proteins (MIA, cystatin M and carbonic anhydrase III) were identified as potentially having a causal effect on eGFR, results which are biologically plausible given their known roles (Supplementary Table 2024). Nevertheless, the precise mechanism through which

1
2
3 these proteins could be exerting an effect on eGFR remains to be elucidated. Discordant
4
5 directions of effect from the observational and the causal estimates (cystatin M, carbonic
6
7 anhydrase III) could be explained due to differences in sample size and characteristics,
8
9 reverse causation or confounding in the case of the observational estimates, or due to
10
11 limitations innate to the MR methods^{86, 87}. A further explanation might be that they are
12
13 different effects: MR examines lifelong exposures to higher/lower protein levels, whereas
14
15 results from observational studies could be reflective of acute effects⁸⁷. The causality inferred
16
17 by MR implies that if the exposure variables are lowered by means of an intervention
18
19 mimicking the effect of the SNP on the exposure, the outcome variable will also change in
20
21 the direction predicted by MR. That is, our results provide suggestive evidence that lifestyle
22
23 or pharmacological interventions designed to improve renal function (with eGFR as a proxy
24
25 thereof) increase plasma testican-2 levels. However, a limitation inherent to MR is that no
26
27 information is offered on the time interval (e.g. during development) or target tissue in which
28
29 such intervention would need to be delivered⁸⁸. Likewise, current knowledge of the biological
30
31 role of the proteins identified is insufficient for our findings to suggest mechanistic insights.
32
33 This represents an opportunity for future research to unravel molecular mechanisms
34
35 potentially underlying findings from MR and the biology of certain proteins in renal function.
36
37

38
39 The strengths of our study include the use of a multiplex proteomics platform and large
40
41 sample size. This is the one of the first study-studies to report eGFR-protein associations
42
43 with replication in independent samples of diverse ancestries, adding to the robustness and
44
45 generalizability of our results. Ours is ~~also~~ the first study to examine causality in proteomic
46
47 associations with eGFR. The MR was conducted with the largest available GWAS results
48
49 from non-overlapping European ancestry populations, thus avoiding issues derived from
50
51 population stratification and sample overlap. We additionally used GWA summary statistics
52
53 from a complementary renal trait (BUN) to improve the specificity of the SNPs used as
54
55 genetic instruments for eGFR. We focused on cis-SNPs in the reverse MR, thus reducing the
56
57 possibility of horizontal pleiotropy. Finally, the conclusions presented here are supported by
58
59 the multiple sensitivity analyses conducted to test the robustness of our methods.
60

1
2
3 Our study also has several limitations. Aptamer-based proteomic methods may be
4 affected by probe cross-reactivity and non-specific binding^{89,33} – yet the we were able to
5 validate the aptamer-based data from 12 proteins were measured using an independent
6 analytical method in a subset of the discovery sample. Moreover, aptamer-based
7 measurements of 50 of our 57 proteins have been validated by mass spectrometry in blood
8 plasma/serum/other biological matrices and solution affinity measurements in multiple
9 independent studies^{27, 38, 39}. This platform does not produce absolute concentrations of
10 plasma proteins, thus limiting the interpretability of the regression coefficients. Future
11 validation studies developing absolute quantitative assays for the detection of testican-2, as
12 well as other proteins identified here, are warranted in order to establish reference ranges
13 and to explore their suitability as prognostic and diagnostic biomarkers in clinical settings.
14 Likewise, proteins undergoing post-transcriptional and post-translational modifications are
15 not covered in this assay (although investigation of the modification of specific reagents for
16 this purpose is underway⁸⁹), so potentially relevant proteins undergoing such molecular
17 modifications may have been overlooked in this study. Our findings are based on cross-
18 sectional data, so future studies examining longitudinal changes in these proteins would be
19 of interest. We examined linear associations with plasma proteomic levels, although non-
20 linear trajectories might also exist for some proteins. Likewise, the age interaction effects
21 observed warrant further investigation in future longitudinal studies. We avoided weak
22 instrument bias in MR by selecting genome-wide significant variants, but cannot discount the
23 possibility of having incurred selection bias in the case of the SNP-protein data. The sample
24 size in which genetic associations with protein levels were calculated was significantly
25 smaller than the sample used to identify genetic associations with eGFR, which likely
26 resulted in differences in power. Moreover, despite the multi-ethnic nature of our study, Asian
27 and African ancestries were not represented, and the generalizability of our results may not
28 extend to these populations. Finally, a follow-up of the findings of our study in appropriate
29 experimental models would provide additional evidence on the inferred causal associations
30 reported here and help to unravel the molecular mechanisms underlying our findings.
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

1
2
3 In summary, we present the largest study to date examining the plasma proteomics of
4 renal function. We use a trans-ethnic discovery-replication setting, and are the first to assess
5 the causality underlying our associations. We identified multiple well-known markers of
6 kidney function, and discovered two proteins not previously known to be associated with
7 eGFR. Our findings with regard to causal relationships represent an important stepping-stone
8 in the establishment of testican-2 as a clinically relevant physiological marker of kidney
9 disease progression, and provide suggestive evidence for a number of proteins causally
10 affecting eGFR warranting further investigation as potential therapeutic targets. Our results
11 may serve as the starting point for future work on the translational role of eGFR-associated
12 proteins as diagnostic or prognostic biomarkers of disease, potentially druggable targets, as
13 well as for research on mechanistic insights at the tissue and single cell levels.

27 **Author contributions**

28
29
30 Conceived and designed the study: P.R.M.G., R.W., J.W., C.G., M.W.

31
32 Performed experiments: D.O.M.K., J.G., Ag.P.

33
34 Analyzed data or interpreted results: P.R.M.G., R.W., Q.G., S.B.Z, J.M.E., Xi.X., F.J.C., J.G.,
35
36 M.P., M.T., M.W.

37
38 Contributed reagents/materials/analysis tools: Q.G., S.B.Z., J.M.E., Xi.X., F.J.C., J.Do., P.S.,
39
40 M.A.E., S.S., A.P., A.F., D.O.M.K., J.W., J.Da., E.D.A., W.-H.-O., N.-A.-W., D.-J.-R., Ag.P.,
41
42 J.G., W.K., K.H., C.J., A.K., A.B., M.P., S.M.H., K.S., C.G., M.T., H.K.T.S.**, A.T., M.W.

43
44 Wrote the paper: P.R.M.G., M.W.

45
46 Edited and improved clarity of the manuscript: R.W., Q.G., S.B.Z., J.M.E., Xi.X., F.J.C., J.Do.,
47
48 P.S., J.N., Ag.P., J.G., W.K., A.K., S.M.H., M.T., A.T., M.W.

49
50 All authors discussed the results, and reviewed and approved the final manuscript.

51 **Acknowledgements**

52
53 We are grateful to all study participants of KORA, HUNT, INTERVAL and QMDiab for their
54
55 invaluable contributions to these studies, as well as all members of field staff conducting the
56
57 studies. An early version of the abstract was sent to the 53rd European Society of Human
58
59
60

1
2
3 [Genetics \(ESHG\) Conference, and is listed online as P03.24.C at the ESHG 2020.2 Abstract](#)
4 [Library \(https://2020.eshg.org/index.php/abstract-library/\)](#)
5
6

7 **Ethics approval and informed consent**

8
9 All participants provided written consent, and details on institutions providing ethics approval
10 from each study are described next:

11
12
13 KORA: The KORA cohort ethical approval was granted by the ethics committee of the
14 Bavarian Medical Association (REC reference number F4: #06068)

15
16
17 HUNT: Regiona Committee for Medical Research Ethics (REK) and Data Inspectorate

18
19
20 INTERVAL: National Research Ethics Service Committee East of England - Cambridge East
21 (Research Ethics Committee (REC) reference 11/EE/0538)

22
23
24 QMDiab: Institutional Review Boards of HMC and WCM-Q under research protocol number
25 11131/11)
26

27 **Disclosures**

28 **Competing interests**

29
30
31 The authors declare the following competing interests: C.J. has received personal fees for
32 research consultancy work from Pfizer and Bayer outside of the submitted work, and J.D. sits
33 on the International Cardiovascular and Metabolic Advisory Board for Novartis (since 2010);
34 the Steering Committee of UK Biobank (since 2011); the MRC International Advisory Group
35 (ING) member, London (since 2013); the MRC High Throughput Science 'Omics Panel
36 Member, London (since 2013); the Scientific Advisory Committee for Sanofi (since 2013); the
37 International Cardiovascular and Metabolism Research and Development Portfolio
38 Committee for Novartis; and the Astra Zeneca Genomics Advisory Board (2018). Other
39 authors declare no competing interests.
40
41
42
43
44
45
46
47
48
49

50 **Funding**

51
52
53 The KORA study was initiated and financed by the Helmholtz Zentrum München – German
54 Research Center for Environmental Health, which is funded by the German Federal Ministry
55 of Education and Research (BMBF) and by the State of Bavaria. This work was also
56 supported by the Biomedical Research Program at Weill Cornell Medicine in Qatar, a
57
58
59
60

1
2
3 program funded by the Qatar Foundation. K.S. is supported by Qatar National Research
4 Fund (QNRF) grant no. NPRPC11-0115-180010. The Nord-Trøndelag Health Study (The
5 HUNT Study) is a collaboration between HUNT Research Centre (Faculty of Medicine,
6 Norwegian University of Science and Technology NTNU), Nord-Trøndelag County Council,
7 Central Norway Health Authority, and the Norwegian Institute of Public Health. The HUNT
8 part of the project re-used protein data that was originally analysed and paid for by
9 Somalogic Inc, CO, USA. Somalogic had no role in the design and conduct of the study;
10 collection of phenotypic data, statistical analysis, and interpretation of the data; preparation,
11 review, or approval of the manuscript; and decision to submit the manuscript for publication.
12 Professor John Danesh is funded by the National Institute for Health Research [Senior
13 Investigator Award]. The views expressed are those of the authors and not necessarily those
14 of the NHS, the NIHR or the Department of Health and Social Care.
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

[RNA-sequencing experiments and kidney gene expression studies were supported by British Heart Foundation project grants \[PG/17/35/33001 and PG/19/16/34270\] and Kidney Research UK grants \[RP_017_20180302 and RP_013_20190305\] to M.T.](#)

Supplemental material - Table of Contents

Supplementary File 1. Supplemental Notes, Figure and Table legends (.docx)

Supplementary File 2. Supplementary Tables (.xlsx)

Supplementary File 3. Supplementary Figures (.pdf)

Supplementary File 4. STROBE Checklist (.pdf)

References

1. Jiang, X, Eales, JM, Scannali, D, Nazgiewicz, A, Prestes, P, Maier, M, et al.: Hypertension and renin-angiotensin system blockers are not associated with expression of angiotensin-converting enzyme 2 (ACE2) in the kidney. *European Heart Journal*, 2020.
2. Eckardt, KU, Coresh, J, Devuyst, O, Johnson, RJ, Kottgen, A, Levey, AS, et al.: Evolving importance of kidney disease: from subspecialty to global health burden. *Lancet (London, England)*, 382: 158-169, 2013.
3. Levey, AS, Inker, LA, Coresh, J: GFR estimation: from physiology to public health. *American journal of kidney diseases : the official journal of the National Kidney Foundation*, 63: 820-834, 2014.
4. Levey, AS, Stevens, LA, Schmid, CH, Zhang, YL, Castro, AF, 3rd, Feldman, HI, et al.: A new equation to estimate glomerular filtration rate. *Annals of internal medicine*, 150: 604-612, 2009.
5. Hill, NR, Fatoba, ST, Oke, JL, Hirst, JA, O'Callaghan, CA, Lasserson, DS, et al.: Global Prevalence of Chronic Kidney Disease - A Systematic Review and Meta-Analysis. *PloS one*, 11: e0158765, 2016.
6. Levey, AS, Atkins, R, Coresh, J, Cohen, EP, Collins, AJ, Eckardt, KU, et al.: Chronic kidney disease as a global public health problem: approaches and initiatives - a position statement from Kidney Disease Improving Global Outcomes. *Kidney Int*, 72: 247-259, 2007.

- 1
2
3 7. Mischak, H, Delles, C, Vlahou, A, Vanholder, R: Proteomic biomarkers in kidney disease:
4
5 issues in development and implementation. *Nature reviews Nephrology*, 11: 221-232,
6
7 2015.
8
- 9 8. Sanchez-Nino, MD, Sanz, AB, Ramos, AM, Fernandez-Fernandez, B, Ortiz, A: Clinical
10
11 proteomics in kidney disease as an exponential technology: heading towards the
12
13 disruptive phase. *Clinical kidney journal*, 10: 188-191, 2017.
14
- 15 9. Good, DM, Zurbig, P, Argiles, A, Bauer, HW, Behrens, G, Coon, JJ, et al.: Naturally
16
17 occurring human urinary peptides for use in diagnosis of chronic kidney disease.
18
19 *Molecular & cellular proteomics : MCP*, 9: 2424-2437, 2010.
20
- 21 10. Gold, L, Ayers, D, Bertino, J, Bock, C, Bock, A, Brody, EN, et al.: Aptamer-based
22
23 multiplexed proteomic technology for biomarker discovery. *PloS one*, 5: e15004,
24
25 2010.
26
27
- 28 11. Ngo, D, Sinha, S, Shen, D, Kuhn, EW, Keyes, MJ, Shi, X, et al.: Aptamer-Based
29
30 Proteomic Profiling Reveals Novel Candidate Biomarkers and Pathways in
31
32 Cardiovascular Disease. *Circulation*, 134: 270-285, 2016.
33
34
- 35 12. Ganz, P, Heidecker, B, Hveem, K, Jonasson, C, Kato, S, Segal, MR, et al.: Development
36
37 and Validation of a Protein-Based Risk Score for Cardiovascular Outcomes Among
38
39 Patients With Stable Coronary Heart Disease. *Jama*, 315: 2532-2541, 2016.
40
- 41 13. Inker, LA, Schmid, CH, Tighiouart, H, Eckfeldt, JH, Feldman, HI, Greene, T, et al.:
42
43 Estimating glomerular filtration rate from serum creatinine and cystatin C. *The New*
44
45 *England journal of medicine*, 367: 20-29, 2012.
46
- 47 14. Pena, MJ, Mischak, H, Heerspink, HJ: Proteomics for prediction of disease progression
48
49 and response to therapy in diabetic kidney disease. *Diabetologia*, 59: 1819-1831,
50
51 2016.
52
- 53 15. Gajjala, PR, Jankowski, V, Heinze, G, Bilo, G, Zanchetti, A, Noels, H, et al.: Proteomic-
54
55 Biostatistic Integrated Approach for Finding the Underlying Molecular Determinants of
56
57 Hypertension in Human Plasma. *Hypertension (Dallas, Tex : 1979)*, 70: 412-419,
58
59 2017.
60

16. Carlsson, AC, Ingelsson, E, Sundstrom, J, Carrero, JJ, Gustafsson, S, Feldreich, T, et al.: Use of Proteomics To Investigate Kidney Function Decline over 5 Years. *Clinical journal of the American Society of Nephrology : CJASN*, 12: 1226-1235, 2017.
17. Christensson, A, Ash, JA, DeLisle, RK, Gaspar, FW, Ostroff, R, Grubb, A, et al.: The Impact of the Glomerular Filtration Rate on the Human Plasma Proteome. *Proteomics Clinical applications*, 2017.
18. Ngo, D, Wen, D, Gao, Y, Keyes, MJ, Drury, ER, Katz, DH, et al.: Circulating testican-2 is a podocyte-derived marker of kidney health. *Proceedings of the National Academy of Sciences of the United States of America*, 117: 25026-25035, 2020.
19. Pierce, BL, Burgess, S: Efficient design for Mendelian randomization studies: subsample and 2-sample instrumental variable estimators. *American journal of epidemiology*, 178: 1177-1184, 2013.
20. Sekula, P, Del Greco M, F, Pattaro, C, Köttgen, A: Mendelian Randomization as an Approach to Assess Causality Using Observational Data. *Journal of the American Society of Nephrology : JASN*, 27: 3253-3265, 2016.
21. Holle, R, Happich, M, Lowel, H, Wichmann, HE: KORA--a research platform for population based health research. *Gesundheitswesen (Bundesverband der Ärzte des Öffentlichen Gesundheitsdienstes (Germany))*, 67 Suppl 1: S19-25, 2005.
22. Wichmann, HE, Gieger, C, Illig, T: KORA-gen--resource for population genetics, controls and a broad spectrum of disease phenotypes. *Gesundheitswesen (Bundesverband der Ärzte des Öffentlichen Gesundheitsdienstes (Germany))*, 67 Suppl 1: S26-30, 2005.
23. Krokstad, S, Langhammer, A, Hveem, K, Holmen, T, Midthjell, K, Stene, T, et al.: Cohort Profile: The HUNT Study, Norway. *International journal of epidemiology*, 42: 968-977, 2012.
24. Moore, C, Bolton, T, Walker, M, Kaptoge, S, Allen, D, Daynes, M, et al.: Recruitment and representativeness of blood donors in the INTERVAL randomised trial assessing varying inter-donation intervals. *Trials*, 17: 458, 2016.

- 1
2
3 25. Mook-Kanamori, DO, Selim, MM, Takiddin, AH, Al-Homsi, H, Al-Mahmoud, KA, Al-
4
5 Obaidli, A, et al.: 1,5-Anhydroglucitol in saliva is a noninvasive marker of short-term
6
7 glycemic control. *The Journal of clinical endocrinology and metabolism*, 99: E479-
8
9 483, 2014.
- 10
11 26. Suhre, K, Arnold, M: Connecting genetic risk to disease end points through the human
12
13 blood plasma proteome. 8: 14357, 2017.
- 14
15 27. Sun, BB, Maranville, JC, Peters, JE, Stacey, D, Staley, JR, Blackshaw, J, et al.: Genomic
16
17 atlas of the human plasma proteome. *Nature*, 558: 73-79, 2018.
- 18
19 28. Naylor, M, Short, MI, Rasheed, H, Lin, H, Jonasson, C, Yang, Q, et al.: Aptamer-Based
20
21 Proteomic Platform Identifies Novel Protein Predictors of Incident Heart Failure and
22
23 Echocardiographic Traits. *Circulation Heart failure*, 13: e006749, 2020.
- 24
25 29. Gold, L, Walker, JJ, Wilcox, SK, Williams, S: Advances in human proteomics at high
26
27 scale with the SOMAscan proteomics platform. *New biotechnology*, 29: 543-549,
28
29 2012.
- 30
31 30. Rohloff, JC, Gelinas, AD, Jarvis, TC, Ochsner, UA, Schneider, DJ, Gold, L, et al.: Nucleic
32
33 Acid Ligands With Protein-like Side Chains: Modified Aptamers and Their Use as
34
35 Diagnostic and Therapeutic Agents. *Molecular therapy Nucleic acids*, 3: e201, 2014.
- 36
37 31. Kraemer, S, Vaught, JD, Bock, C, Gold, L, Katilius, E, Keeney, TR, et al.: From
38
39 SOMAmer-based biomarker discovery to diagnostic and clinical applications: a
40
41 SOMAmer-based, streamlined multiplex proteomic assay. *PloS one*, 6: e26332-
42
43 e26332, 2011.
- 44
45 32. Pattaro, C, Riegler, P, Stifter, G, Modenese, M, Minelli, C, Pramstaller, PP: Estimating
46
47 the glomerular filtration rate in the general population using different equations:
48
49 effects on classification and association. *Nephron Clinical practice*, 123: 102-111,
50
51 2013.
- 52
53 33. Coresh, J, Turin, TC, Matsushita, K, Sang, Y, Ballew, SH, Appel, LJ, et al.: Decline in
54
55 estimated glomerular filtration rate and subsequent risk of end-stage renal disease
56
57 and mortality. *Jama*, 311: 2518-2531, 2014.
- 58
59
60

- 1
2
3 34. K/DOQI clinical practice guidelines for chronic kidney disease: evaluation, classification,
4 and stratification. *American journal of kidney diseases : the official journal of the*
5
6 *National Kidney Foundation*, 39: S1-266, 2002.
7
8
9 35. R Core Team: R: A language and environment for statistical computing. Vienna, Austria,
10 R Foundation for Statistical Computing, 2019.
11
12
13 36. Zaykin, DV: Optimally weighted Z-test is a powerful method for combining probabilities in
14 meta-analysis. *Journal of evolutionary biology*, 24: 1836-1841, 2011.
15
16
17 37. Huth, C, von Toerne, C, Schederecker, F, de Las Heras Gala, T, Herder, C, Kronenberg,
18 F, et al.: Protein markers and risk of type 2 diabetes and prediabetes: a targeted
19 proteomics approach in the KORA F4/FF4 study. *European journal of epidemiology*,
20 34: 409-422, 2019.
21
22
23 38. Williams, SA, Kivimaki, M, Langenberg, C, Hingorani, AD, Casas, JP, Bouchard, C, et al.:
24 Plasma protein patterns as comprehensive indicators of health. *Nature medicine*, 25:
25 1851-1857, 2019.
26
27
28 39. Emilsson, V, Ilkov, M, Lamb, JR, Finkel, N, Gudmundsson, EF, Pitts, R, et al.: Co-
29 regulatory networks of human serum proteins link genetics to disease. *Science (New*
30 *York, NY)*, 361: 769-773, 2018.
31
32
33 40. Kyritsis, KA, Wang, B, Sullivan, J, Lyne, R, Micklem, G: InterMineR: an R package for
34 InterMine databases. *Bioinformatics*, 2019.
35
36
37 41. Huang, DW, Sherman, BT, Lempicki, RA: Systematic and integrative analysis of large
38 gene lists using DAVID bioinformatics resources. *Nature Protocols*, 4: 44-57, 2009.
39
40
41 42. Durinck, S, Spellman, PT, Birney, E, Huber, W: Mapping identifiers for the integration of
42 genomic datasets with the R/Bioconductor package biomaRt. *Nat Protoc*, 4: 1184-
43 1191, 2009.
44
45
46 43. Schmidt, T, Samaras, P, Frejno, M, Gessulat, S, Barnert, M, Kienegger, H, et al.:
47 ProteomicsDB. *Nucleic acids research*, 46: D1271-D1281, 2017.
48
49
50
51
52
53
54
55
56
57
58
59
60

- 1
2
3 44. Samaras, P, Schmidt, T, Frejno, M, Gessulat, S, Reinecke, M, Jarzab, A, et al.:
4
5 ProteomicsDB: a multi-omics and multi-organism resource for life science research.
6
7 *Nucleic acids research*, 48: D1153-D1163, 2019.
8
- 9 45. Consortium, GT: The Genotype-Tissue Expression (GTEx) project. *Nature genetics*, 45:
10
11 580-585, 2013.
12
- 13 46. Szklarczyk, D, Gable, AL, Lyon, D, Junge, A, Wyder, S, Huerta-Cepas, J, et al.: STRING
14
15 v11: protein-protein association networks with increased coverage, supporting
16
17 functional discovery in genome-wide experimental datasets. *Nucleic acids research*,
18
19 47: D607-d613, 2019.
20
- 21 47. Davey Smith, G, Hemani, G: Mendelian randomization: genetic anchors for causal
22
23 inference in epidemiological studies. *Human molecular genetics*, 23: R89-98, 2014.
24
- 25 48. Smith, GD, Ebrahim, S: 'Mendelian randomization': can genetic epidemiology contribute
26
27 to understanding environmental determinants of disease? *International journal of*
28
29 *epidemiology*, 32: 1-22, 2003.
30
- 31 49. Lawlor, DA, Harbord, RM, Sterne, JA, Timpson, N, Davey Smith, G: Mendelian
32
33 randomization: using genes as instruments for making causal inferences in
34
35 epidemiology. *Statistics in medicine*, 27: 1133-1163, 2008.
36
37
- 38 50. Wuttke, M, Li, Y, Li, M, Sieber, KB, Feitosa, MF, Gorski, M, et al.: A catalog of genetic
39
40 loci associated with kidney function from analyses of a million individuals. *Nature*
41
42 *genetics*, 51: 957-972, 2019.
43
- 44 51. Hemani, G, Zheng, J, Elsworth, B, Wade, KH, Haberland, V, Baird, D, et al.: The MR-
45
46 Base platform supports systematic causal inference across the human phenome.
47
48 *eLife*, 7, 2018.
49
- 50 52. Buniello, A, MacArthur, JAL, Cerezo, M, Harris, LW, Hayhurst, J, Malangone, C, et al.:
51
52 The NHGRI-EBI GWAS Catalog of published genome-wide association studies,
53
54 targeted arrays and summary statistics 2019. *Nucleic acids research*, 47: D1005-
55
56 d1012, 2019.
57
58
59
60

- 1
2
3 53. Burgess, S, Butterworth, A, Thompson, SG: Mendelian randomization analysis with
4 multiple genetic variants using summarized data. *Genetic epidemiology*, 37: 658-665,
5 2013.
6
7
8
9 54. Teumer, A: Common Methods for Performing Mendelian Randomization. *Frontiers in*
10 *cardiovascular medicine*, 5: 51, 2018.
11
12
13 55. Burgess, S, Bowden, J, Fall, T, Ingelsson, E, Thompson, SG: Sensitivity Analyses for
14 Robust Causal Inference from Mendelian Randomization Analyses with Multiple
15 Genetic Variants. *Epidemiology*, 28: 30-42, 2017.
16
17
18
19 56. Bowden, J, Davey Smith, G, Burgess, S: Mendelian randomization with invalid
20 instruments: effect estimation and bias detection through Egger regression.
21 *International journal of epidemiology*, 44: 512-525, 2015.
22
23
24
25 57. Bowden, J, Davey Smith, G, Haycock, PC, Burgess, S: Consistent Estimation in
26 Mendelian Randomization with Some Invalid Instruments Using a Weighted Median
27 Estimator. *Genetic epidemiology*, 40: 304-314, 2016.
28
29
30
31 58. Hartwig, FP, Davey Smith, G, Bowden, J: Robust inference in summary data Mendelian
32 randomization via the zero modal pleiotropy assumption. *International journal of*
33 *epidemiology*, 46: 1985-1998, 2017.
34
35
36
37 59. Verbanck, M, Chen, C-Y, Neale, B, Do, R: Detection of widespread horizontal pleiotropy
38 in causal relationships inferred from Mendelian randomization between complex traits
39 and diseases. *Nature genetics*, 50: 693-698, 2018.
40
41
42
43 60. Ju, W, Nair, V, Smith, S, Zhu, L, Shedden, K, Song, PXK, et al.: Tissue transcriptome-
44 driven identification of epidermal growth factor as a chronic kidney disease biomarker.
45 *Science translational medicine*, 7: 316ra193, 2015.
46
47
48
49 61. Martini, S, Eichinger, F, Nair, V, Kretzler, M: Defining human diabetic nephropathy on the
50 molecular level: Integration of transcriptomic profiles with biological knowledge.
51 *Reviews in Endocrine and Metabolic Disorders*, 9: 267-274, 2008.
52
53
54
55 62. Sampson, MG, Robertson, CC, Martini, S, Mariani, LH, Lemley, KV, Gillies, CE, et al.:
56 Integrative Genomics Identifies Novel Associations with APOL1 Risk Genotypes in
57
58
59
60

- 1
2
3 Black NEPTUNE Subjects. *Journal of the American Society of Nephrology : JASN*,
4
5 27: 814-823, 2016.
6
- 7 63. Reich, HN, Tritchler, D, Cattran, DC, Herzenberg, AM, Eichinger, F, Boucherot, A, et al.:
8
9 A molecular signature of proteinuria in glomerulonephritis. *PLoS one*, 5: e13451,
10
11 2010.
12
- 13 64. Rodwell, GE, Sonu, R, Zahn, JM, Lund, J, Wilhelmy, J, Wang, L, et al.: A transcriptional
14
15 profile of aging in the human kidney. *PLoS biology*, 2: e427, 2004.
16
- 17 65. Balduzzi, S, Rücker, G, Schwarzer, G: How to perform a meta-analysis with R: a practical
18
19 tutorial. *Evidence Based Mental Health*, 22: 153-160, 2019.
20
- 21 66. Xu, X, Eales, JM, Akbarov, A, Guo, H, Becker, L, Talavera, D, et al.: Molecular insights
22
23 into genome-wide association studies of chronic kidney disease-defining traits.
24
25 *Nature communications*, 9: 4800, 2018.
26
- 27 67. Rowland, J, Akbarov, A, Eales, J, Xu, X, Dormer, JP, Guo, H, et al.: Uncovering genetic
28
29 mechanisms of kidney aging through transcriptomics, genomics, and epigenomics.
30
31 *Kidney Int*, 95: 624-635, 2019.
32
- 33 68. Richmond, RC, Davey Smith, G: Commentary: Orienting causal relationships between
34
35 two phenotypes using bidirectional Mendelian randomization. *International journal of*
36
37 *epidemiology*, 48: 907-911, 2019.
38
- 39 69. Lehallier, B, Gate, D, Schaum, N, Nanasi, T, Lee, SE, Yousef, H, et al.: Undulating
40
41 changes in human plasma proteome profiles across the lifespan. *Nature medicine*,
42
43 25: 1843-1850, 2019.
44
- 45 70. Jovanovic, D, Krstivojevic, P, Obradovic, I, Durdevic, V, Dukanovic, L: Serum cystatin C
46
47 and beta2-microglobulin as markers of glomerular filtration rate. *Renal failure*, 25:
48
49 123-133, 2003.
50
- 51 71. Niewczas, MA, Pavkov, ME, Skupien, J, Smiles, A, Md Dom, ZI, Wilson, JM, et al.: A
52
53 signature of circulating inflammatory proteins and development of end-stage renal
54
55 disease in diabetes. *Nature medicine*, 25: 805-813, 2019.
56
57
58
59
60

- 1
2
3 72. Prunotto, M, Farina, A, Lane, L, Pernin, A, Schifferli, J, Hochstrasser, DF, et al.:
4
5 Proteomic analysis of podocyte exosome-enriched fraction from normal human urine.
6
7 *Journal of proteomics*, 82: 193-229, 2013.
8
- 9 73. Jia, L, Zhang, L, Shao, C, Song, E, Sun, W, Li, M, et al.: An attempt to understand
10
11 kidney's protein handling function by comparing plasma and urine proteomes. *PLoS*
12
13 *one*, 4: e5146, 2009.
14
- 15 74. Schenk, S, Schoenhals, GJ, de Souza, G, Mann, M: A high confidence, manually
16
17 validated human blood plasma protein reference set. *BMC medical genomics*, 1: 41,
18
19 2008.
20
- 21 75. Carlsson, AC, Sundström, J, Carrero, JJ, Gustafsson, S, Stenemo, M, Larsson, A, et al.:
22
23 Use of a proximity extension assay proteomics chip to discover new biomarkers
24
25 associated with albuminuria. *European journal of preventive cardiology*, 24: 340-348,
26
27 2016.
28
- 29 76. Mohammadi-Shemirani, P, Sjaarda, J, Gerstein, HC, Treleaven, DJ, Walsh, M, Mann, JF,
30
31 et al.: A Mendelian Randomization-Based Approach to Identify Early and Sensitive
32
33 Diagnostic Biomarkers of Disease. *Clinical chemistry*, 65: 427-436, 2019.
34
35
- 36 77. Clark, CJ, Sage, EH: A prototypic matricellular protein in the tumor microenvironment--
37
38 where there's SPARC, there's fire. *Journal of cellular biochemistry*, 104: 721-732,
39
40 2008.
41
- 42 78. Feng, D, Ngov, C, Henley, N, Boufaied, N, Gerarduzzi, C: Characterization of
43
44 Matricellular Protein Expression Signatures in Mechanistically Diverse Mouse Models
45
46 of Kidney Injury. *Scientific reports*, 9: 16736, 2019.
47
48
- 49 79. Marimuthu, A, O'Meally, RN, Chaerkady, R, Subbannayya, Y, Nanjappa, V, Kumar, P, et
50
51 al.: A comprehensive map of the human urinary proteome. *Journal of proteome*
52
53 *research*, 10: 2734-2743, 2011.
54
- 55 80. Francki, A, Sage, EH: SPARC and the kidney glomerulus: matricellular proteins exhibit
56
57 diverse functions under normal and pathological conditions. *Trends in cardiovascular*
58
59 *medicine*, 11: 32-37, 2001.
60

- 1
2
3 81. Woroniecka, KI, Park, AS, Mohtat, D, Thomas, DB, Pullman, JM, Susztak, K:
4
5 Transcriptome analysis of human diabetic kidney disease. *Diabetes*, 60: 2354-2369,
6
7 2011.
8
9 82. Lindenmeyer, MT, Eichinger, F, Sen, K, Anders, HJ, Edenhofer, I, Mattinzoli, D, et al.:
10
11 Systematic analysis of a novel human renal glomerulus-enriched gene expression
12
13 dataset. *PloS one*, 5: e11545, 2010.
14
15 83. Nystrom, J, Fierlbeck, W, Granqvist, A, Kulak, SC, Ballermann, BJ: A human glomerular
16
17 SAGE transcriptome database. *BMC nephrology*, 10: 13, 2009.
18
19 84. Ju, W, Greene, CS, Eichinger, F, Nair, V, Hodgins, JB, Bitzer, M, et al.: Defining cell-type
20
21 specificity at the transcriptional level in human disease. *Genome research*, 23: 1862-
22
23 1873, 2013.
24
25 85. Wynn, TA, Ramalingam, TR: Mechanisms of fibrosis: therapeutic translation for fibrotic
26
27 disease. *Nature medicine*, 18: 1028-1040, 2012.
28
29 86. Haycock, PC, Burgess, S, Wade, KH, Bowden, J, Relton, C, Davey Smith, G: Best (but
30
31 oft-forgotten) practices: the design, analysis, and interpretation of Mendelian
32
33 randomization studies. *The American journal of clinical nutrition*, 103: 965-978, 2016.
34
35 87. Davies, NM, Holmes, MV, Davey Smith, G: Reading Mendelian randomisation studies: a
36
37 guide, glossary, and checklist for clinicians. *BMJ (Clinical research ed)*, 362: k601,
38
39 2018.
40
41 88. Bretherick, AD, Canela-Xandri, O, Joshi, PK, Clark, DW, Rawlik, K, Boutin, TS, et al.:
42
43 Linking protein to phenotype with Mendelian Randomization detects 38 proteins with
44
45 causal roles in human diseases and traits. *PLoS genetics*, 16: e1008785, 2020.
46
47 89. Suhre, K, McCarthy, MI, Schwenk, JM: Genetics meets proteomics: perspectives for
48
49 large population-based studies. *Nature reviews Genetics*, 2020.
50
51
52
53
54
55
56
57
58
59
60

Tables

Table 1. Population characteristics of association studies

Trait	KORA	HUNT3	INTERVAL	QMDIAB
N	995	930	623	334
Age (years)	59.31 (7.81)	68.94 (10.29)	47.36 (13.35)	47.10 (12.57)
Male	480 (48.2)	688 (74.0)	343 (55.1)	169 (50.6)
BMI (kg/m ²)	27.78 (4.58)	28.38 (3.97)	27.16 (10.05)	29.66 (5.95)
Smoking	572 (57.5)	699 (75.16)	99 (15.89)	60 (18.0)
Serum creatinine (mg/dL)	0.85 (0.18)	0.92 (0.32)	0.70 (0.14)	0.85 (0.22)
eGFR (mL/min/1.73m ²)	85.98 (14.06)	80.25 (18.75)	108.27 (16.21)	95.87 (27.32)
CKD	38 (3.8)	138 (14.8)	1 (0.16%)	30 (9.0)
UACR (mg/dL) *	5.64 (3.61, 9.94)	NA	NA	NA
MA	58 (5.9)	NA	NA	NA
HDL cholesterol (mg/dL)	57.32 (15.20)	45.12 (11.24)	74.33 (24.42)	47.58 (13.75)
Triglycerides (mg/dL) *	107 (75, 155.5)	141.27 (106.28, 194.85)	132.75 (97.35, 194.70)	169 (99.20, 215.23)
Anti-hyperlipidemic medication use	142 (14.3)	NA	33 (5.30)	NA
Hypertension	397 (39.9)	355 (38.2)	48 (7.70)	103 (30.8)
Diabetes	68 (6.8)	128 (13.8)	2 (0.32)	172 (51.5)

Table legend: Measurement units are shown in parentheses in the trait column, where the absence of units means it is a categorical trait. The mean and (SD) are presented for non-skewed continuous variables, while skewed continuous variables are identified with * and median (1st, 3rd quartile) are presented. Count and (%) are shown for categorical variables. HUNT3: third survey of Nord-Trøndelag Health Study (HUNT3), INTERVAL: INTERVAL Study, QMDIAB: Qatar Metabolomics Study on Diabetes; NA: not available; BMI: body-mass index, eGFR: serum creatinine-based estimated glomerular filtration rate, CKD: chronic kidney disease, UACR: urinary albumin-to-creatinine ratio, MA: microalbuminuria, HDL: high density lipoprotein.

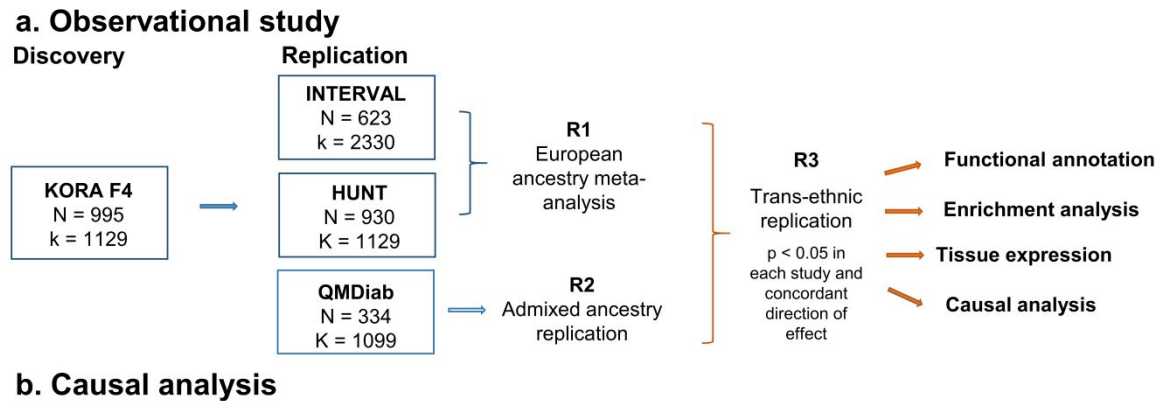
Table 2. Causal estimates across Mendelian randomization methods

Association		IVW	MR-Egger	Weighted median	Weighted mode	Wald ratio
eGFR → SPOCK2	β	5.213	8.61	6.104	6.506	
	CI	[2.767, 7.659]	[2.405, 14.815]	[2.808, 9.401]	[1.189, 11.823]	--
	p	2.95E-05	0.01	2.84E-04	0.021	
	N _{IV}	40	40	40	40	
CA3 → eGFR	β					0.007
	CI	--	--	--	--	[0.003, 0.010]
	p					5.04E-04
	N _{IV}					1
CST6 → eGFR	β					0.007
	CI	--	--	--	--	[0.004, 0.011]
	p					8.41E-05
	N _{IV}					1
MIA → eGFR	β	-0.002	-0.001	-0.002	-0.001	
	CI	[-0.003, -0.001]	[-0.002, 0.000]	[-0.003, -0.001]	[-0.002, -0.001]	--
	p	8.79E-04	0.299	2.00E-04	0.081	
	N _{IV}	3	3	3	3	

Table legend: Results from the forward MR (effect of eGFR on protein levels, i.e. eGFR → protein) are based on the 40 instruments retrieved from Wuttke, et.al. 2019⁵⁰, whereas the reverse MR (protein → eGFR) are based on the one to three instruments retrieved from [the INTERVAL pGWAS reported in Sun, et.al. 2018⁵⁸](#). In bold are significant p values at a Bonferroni-corrected level (0.05/47 for the forward analysis, 0.05/28 for the reverse analysis). *SPOCK2*: testican-2, *CA3*: carbonic anhydrase III, *CST6*: cystatin-M, *MIA*: melanoma-derived growth regulatory protein, IVW: Inverse-variance weighted MR; ~~MR-PRESSO: Mendelian randomization Pleiotropy RESidual Sum and Outlier~~; β : causal estimate, SE: Standard Error, CI: 95% confidence interval of causal estimate; N_{IV}: number of SNPs used as instrumental variables.

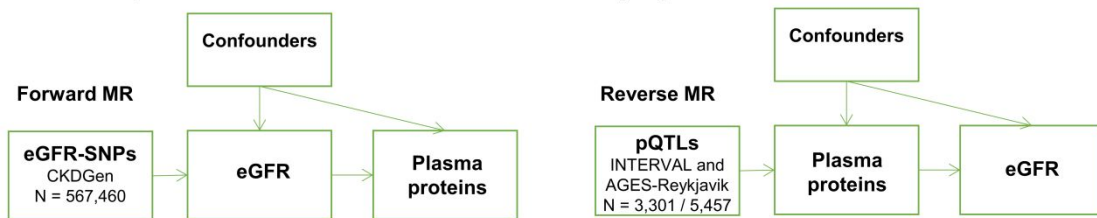
Figures

Figure 1. Study overview



b. Causal analysis

Two-sample bidirectional Mendelian Randomization (MR)



MR analysis flowchart



Figure 2. Cross-sectional effect estimates from the top 10 proteins associated with eGFR

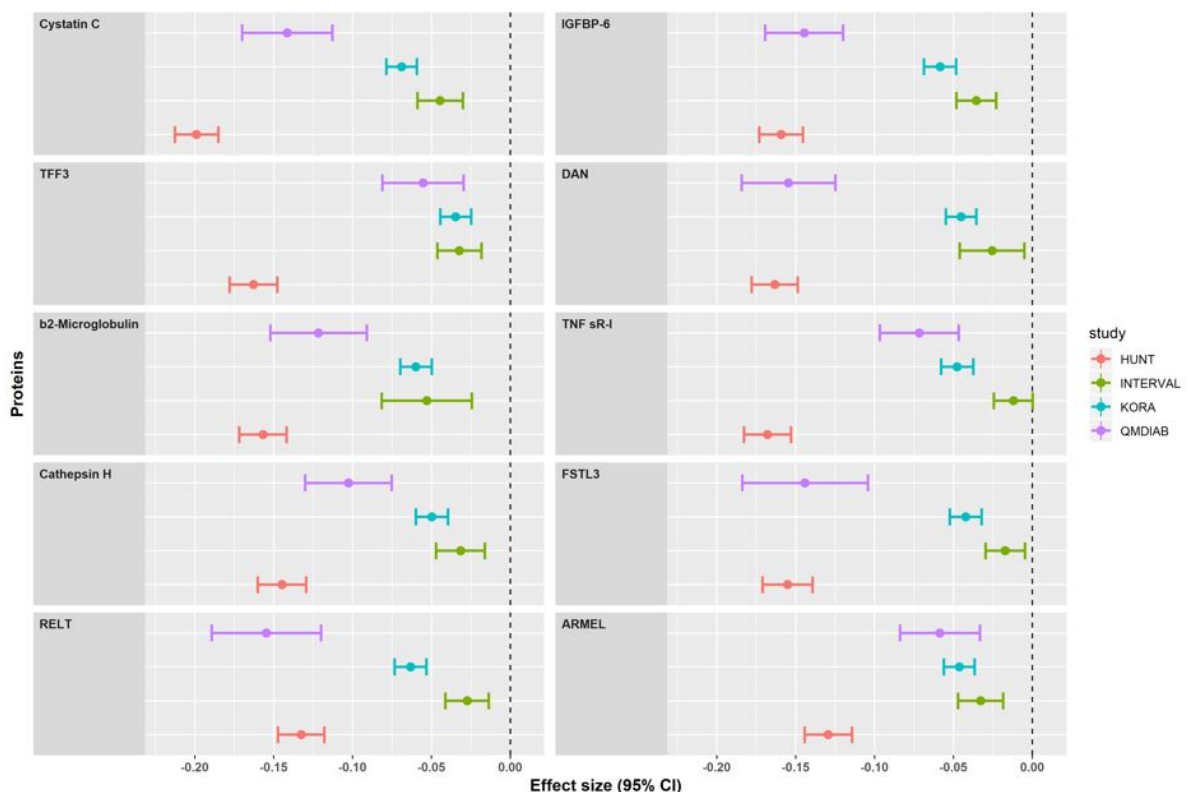


Figure 3. Overlap of associations with eGFR and CKD

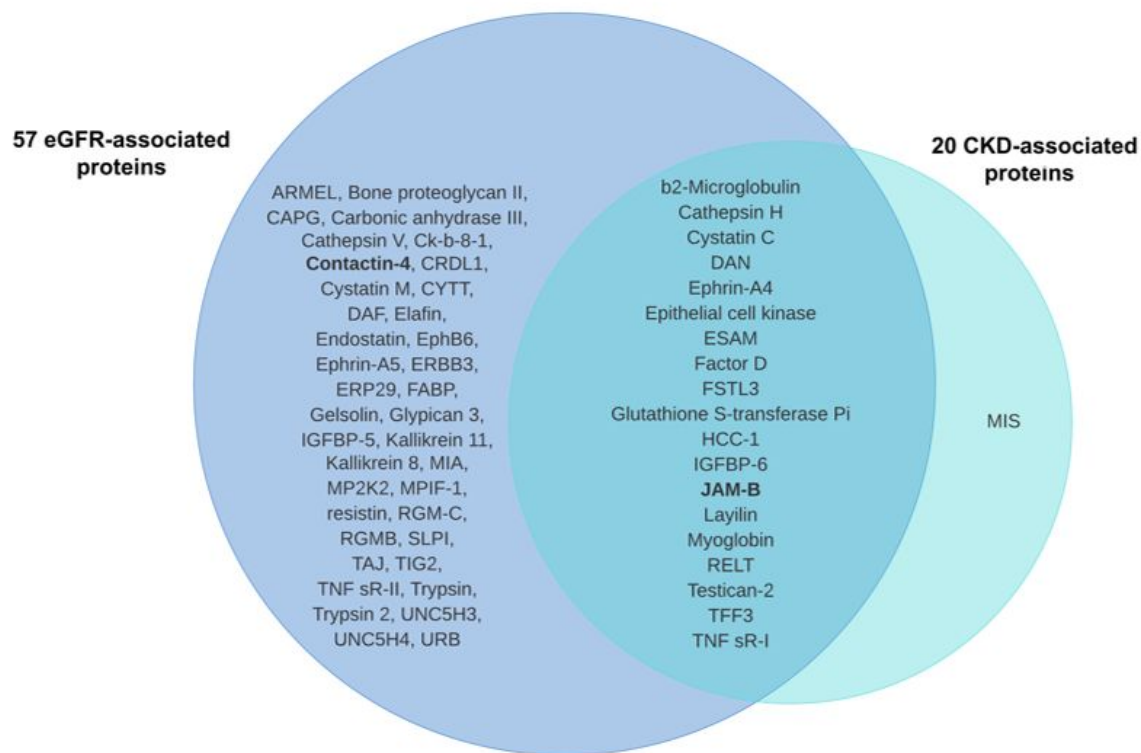


Figure 4. Results from the forward MR estimating the causal effect of eGFR on testican-2

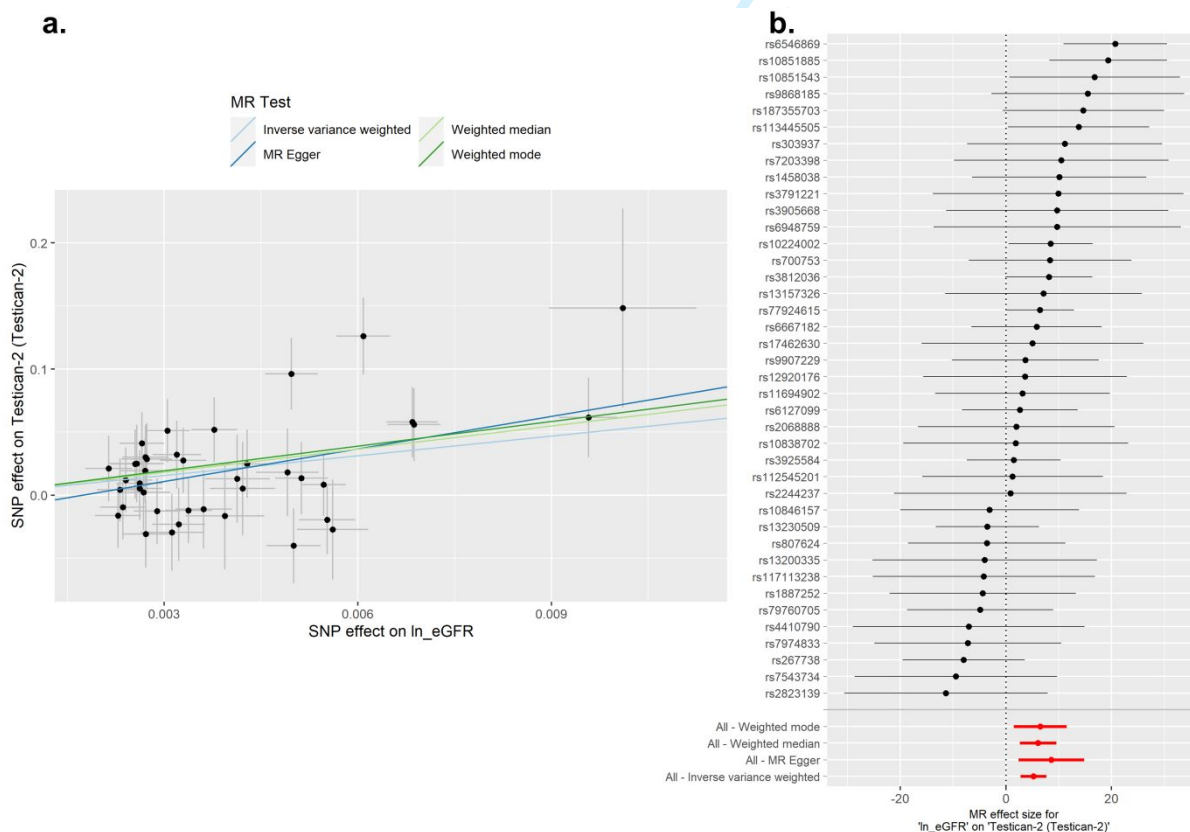


Figure legends

Fig. 1. Panel a. Cross-sectional association study. Data from 995 participants and 1129 proteins from KORA F4 was used in the discovery phase of a proteome wide association study of renal function using confounder-adjusted regression models. The replication studies were INTERVAL, HUNT and QMDiab. Three rounds of replication are shown: R1, replication based on the meta-analysis of p-values from the linear regression results of the studies with European ancestry; R2, replication based on the results of linear regression models performed in the Arab, South Asian and Filipino descent sample QMDiab; R3, identification of proteins consistently associated with eGFR across samples and ethnicities. The set of proteins identified in R3 was then functionally annotated and brought forward to the causal analysis phase. **Panel b. Causal analysis.** Two-sample bidirectional Mendelian randomization using data on participants from EA-studies in the CKDGen Consortium to instrument the forward analysis (eGFR causal to protein level) and data from INTERVAL and AGES-Reykjavik to instrument the reverse analysis (protein level causal to eGFR). [Details on the data processing workflow for Mendelian randomizations are shown.](#)

Fig. 2. Regression coefficient estimates from the top 10 proteins identified in the cross-sectional association trans-ethnic study on eGFR. The x-axis shows the estimates and 95% CI for the regression coefficients (i.e. change in log-transformed eGR per standard deviation increase in protein level), and each panel corresponds to one protein. Estimates are color coded according to the specific study: HUNT3 in red, INTERVAL in green, KORA in blue, and QMDiab in purple. TFF3: Trefoil factor 3, RELT: Tumor necrosis factor receptor superfamily member 19L, IGFBP-6: Insulin-like growth factor-binding protein 6, DAN: Neuroblastoma suppressor of tumorigenicity 1, TNF sR-I: Tumor necrosis factor receptor superfamily member 1A, FSTL3: Follistatin-related protein 3, ARMEL: Cerebral dopamine neurotrophic factor.

Fig. 3. Results from the trans-ethnic discovery-replication observational study. Depicted in the left circle are the 57 proteins associated with eGFR, the continuous measurement of renal function; the 38 eGFR-specific proteins reflect associations along the full range of renal

1
2
3 function, whereas the 19 proteins also associated with CKD reflect a direct association with a
4 clinically relevant low eGFR (<60 mL/min/1.73m²). In bold are the two novel proteins found
5 by this study.
6
7

8
9 **Fig.4. Panel a.** Scatter plot showing the individual genetic effects of the selected
10 instrumental variables on log transformed eGFR (coefficient of the SNP-exposure
11 association) on the x-axis and on testican-2 plasma levels (coefficient of the SNP-outcome)
12 on the y-axis, along with their 95% CI. Each data point corresponds to an individual SNP.
13
14 The lines correspond to the slopes of the different MR methods, which can be interpreted as
15 the change in testican-2 levels per unit increase in log-transformed eGFR, and are color
16 coded as follows: IVW-MR in light blue, MR-Egger in dark blue, weighted median in light
17 green, weighted mode in dark green. **Panel b.** Forest plot showing the individual causal
18 estimates of each of the 40 genetic instruments. The red points show the pooled estimates
19 using all SNPs in the four methods. 95% CI are shown. IVW: inverse-variance-weighted; MR:
20 Mendelian randomization.
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

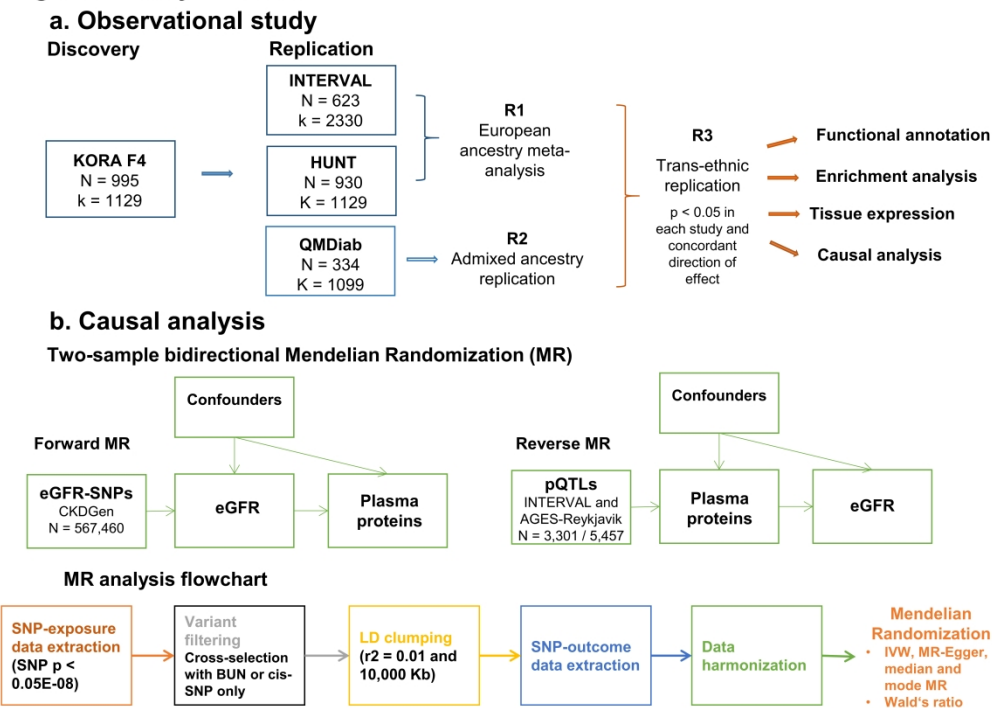
Figure 1. Study overview

Fig. 1. Panel a. Cross-sectional association study. Data from 995 participants and 1129 proteins from KORA F4 was used in the discovery phase of a proteome wide association study of renal function using confounder-adjusted regression models. The replication studies were INTERVAL, HUNT and QMDiab. Three rounds of replication are shown: R1, replication based on the meta-analysis of p-values from the linear regression results of the studies with European ancestry; R2, replication based on the results of linear regression models performed in the Arab, South Asian and Filipino descent sample QMDiab; R3, identification of proteins consistently associated with eGFR across samples and ethnicities. The set of proteins identified in R3 was then functionally annotated and brought forward to the causal analysis phase. Panel b. Causal analysis. Two-sample bidirectional Mendelian randomization using data on participants from EA-studies in the CKDGen Consortium to instrument the forward analysis (eGFR causal to protein level) and data from INTERVAL and AGES-Reykjavik to instrument the reverse analysis (protein level causal to eGFR). Details on the data processing workflow for Mendelian randomizations are shown.

Figure 4. Results from forward MR estimating the causal effect of eGFR on testican-2

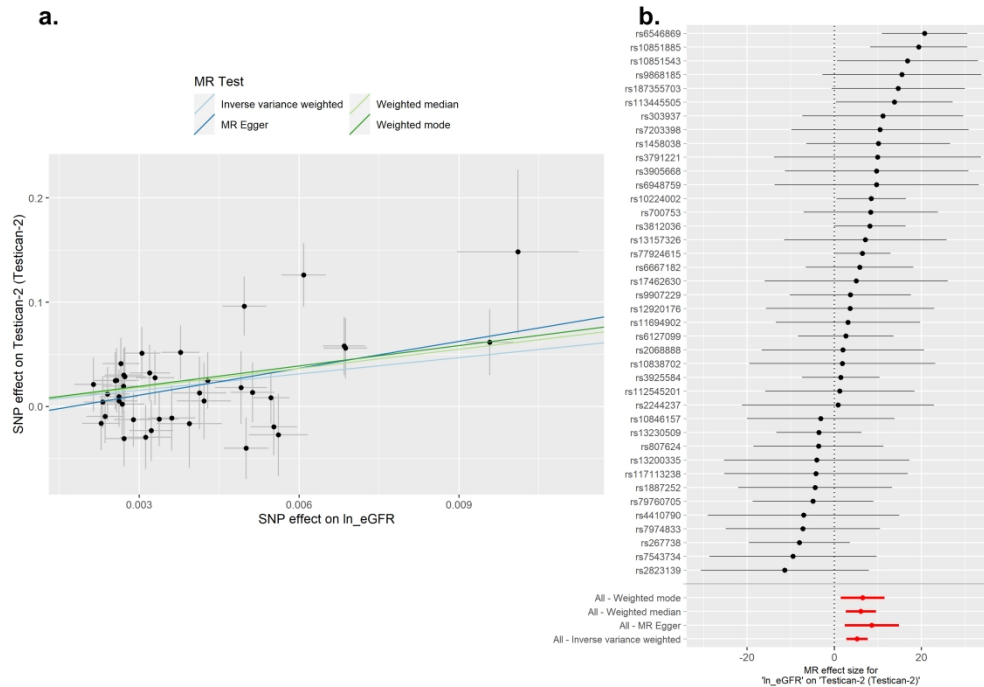


Fig.4. Panel a. Scatter plot showing the individual genetic effects of the selected instrumental variables on log transformed eGFR (coefficient of the SNP-exposure association) on the x-axis and on testican-2 plasma levels (coefficient of the SNP-outcome) on the y-axis, along with their 95% CI. Each data point corresponds to an individual SNP. The lines correspond to the slopes of the different MR methods, which can be interpreted as the change in testican-2 levels per unit increase in log-transformed eGFR, and are color coded as follows: IVW-MR in light blue, MR-Egger in dark blue, weighted median in light green, weighted mode in dark green. Panel b. Forest plot showing the individual causal estimates of each of the 40 genetic instruments. The red points show the pooled estimates using all SNPs in the four methods. 95% CI are shown. IVW: inverse-variance-weighted; MR: Mendelian randomization.

**Plasma proteomics of renal function: a trans-ethnic meta-analysis and Mendelian
randomization study**

Supplemental Material

Supplemental Notes	4
Supplemental Note 1: Data availability	4
Supplemental Note 2: Covariate definition	4
Supplemental Note 3: Interaction analyses	5
Supplemental Note 4: Mendelian randomization analysis	6
Supplemental Note 5: Expression analyses in human kidney samples	10
Legends to Supplementary Tables	13
Supplementary Table 1. Results from cross-sectional analysis of eGFR and CKD in KORA F4.....	13
Supplementary Table 2. Validation of proteomic targets from Emilsson et.al., 2018 ..	13
Supplementary Table 3. Validation of proteomic targets from Sun et.al., 2018	13
Supplementary Table 4. Validation of proteomic targets from Williams et.al., 2019.....	13
Supplementary Table 5. Genetic instrument selection and data harmonization (INTERVAL and AGES- Reykjavik)	14
Supplementary Table 6. Harmonized summary statistics used in MR	14
Supplementary Table 7. Gene-trait information retrieved from GWAS Catalog for pleiotropic SNPs	14
Supplementary Table 8. Sensitivity analyses with CysC-based eGFR, no adjustment for BMI and no adjustment for T2D	15
Supplementary Table 9. Sensitivity analyses (interaction with age, sex and smoking) in KORA F4.....	15
Supplementary Table 10. Results from observational analysis of supplementary renal phenotypes (eGFR decline, log(uACR) and MA) in KORA F4	15
Supplementary Table 11. Replication of cross-sectional eGFR-protein associations ..	15
Supplementary Table 12. Replication of cross-sectional CKD-protein associations ..	16
Supplementary Table 13. Extended annotation file (DAVID)	16
Supplementary Table 14. Gene information from proteins included in MR	16
Supplementary Table 15. Phenotypic variance explained by MR instruments	17

1		
2		
3	Supplementary Table 16. Results from MR (IVW, MBE, weighted median and MR-	
4	Egger)	17
5		
6	Supplementary Table 17. Sensitivity analyses: Heterogeneity (Cochran's Q test)	17
7		
8	Supplementary Table 18. Sensitivity analyses: Pleiotropy in MR-Egger	17
9		
10	Supplementary Table 19. Sensitivity analyses: Leave-one-out analyses	18
11		
12	Supplementary Table 20. Sensitivity analyses: Results from restrictive MR.....	18
13		
14	Supplementary Table 21. Results from correlation analyses between gene	
15	expression and eGFR from Nephroseq datasets	18
16		
17	Supplementary Table 22. Clinical characteristics of studies included in gene	
18	expression analyses	18
19		
20	Supplementary Table 23. Results from multivariate regression analyses on gene	
21	expression, eGFR and histological characteristic scoring from human kidney	
22	resource	19
23		
24	Supplementary Table 24. Description and biological roles of selected proteins.....	19
25		
26	Legends to Supplementary Figures.....	20
27		
28	Supplementary Figure 1. Correlation of serum creatinine variables in KORA F4	20
29		
30	Supplementary Figure 2. Correlation between aptamer-based and other	
31	measurements for proteins in KORA F4.....	20
32		
33	Supplementary Figure 3. Genetic instrument selection and data harmonization.....	20
34		
35	Supplementary Figure 4. Protein overlap in pGWAS datasets used in reverse	
36	direction of MR	20
37		
38	Supplementary Figure 5. Proteins and log(eGFR) distribution in discovery dataset ..	21
39		
40	Supplementary Figure 6. Cross sectional results for eGFR-protein associations	
41	across studies	21
42		
43	Supplementary Figure 7. Proteins and log(eGFR) distribution in discovery dataset	
44	after CKD exclusion	21
45		
46	Supplementary Figure 8. Correlation between Z-values for eGFR-protein associations	
47	across studies	22
48		
49	Supplementary Figure 9. Tissue expression of 57 eGFR-associated proteins	
50	(ProteomeDB)	22
51		
52	Supplementary Figure 10. Tissue expression of 56 eGFR-associated protein coding	
53	genes (ProteomeDB).....	22
54		
55	Supplementary Figure 11. Expression of 56 eGFR-associated protein coding genes	
56	across tissues (GTEx)	23
57		
58	Supplementary Figure 12. Protein-protein interaction network of 57 replicated eGFR-	
59	associated proteins.....	23
60		
	Supplementary Figures 13-19. Forward MR results for effects of eGFR on proteins.	23

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

Supplementary Figure 20. Reverse MR analysis for MIA-eGFR23
Supplementary Figure 21. SPOCK2 gene expression in renal tissue from 26 CKD patients.....24

References 25

For Peer Review

Supplemental Notes

Supplemental Note 1: Data availability

All data generated during this study are included in this published article and its supplementary information files. The analysis code in R is available on request. The informed consent given by the study participants does not cover posting of participant level phenotype data in public databases. Pre-existing data access policies for each of the four studies state research data requests can be submitted to each steering committee. Study-specific details regarding such requests are described in the next paragraphs:

KORA data are available upon request from KORA Project Application Self-Service Tool (<https://epi.helmholtz-muenchen.de/>); data requests can be submitted online and are subject to approval by the KORA Board. Data of the QMDiab study will be shared with researchers whose requests have been approved by the KORA Board.

The Nord-Trøndelag Health Study (HUNT) holds comprehensive data from more than 145,000 persons. HUNT Research Centre has been given concession to store and handle these data by the Norwegian Data Inspectorate. The key identification in the data base is the personal identification number given to all Norwegians at birth or immigration, whilst de-identified data are sent to researchers. Due to confidentiality HUNT Research Centre wants to limit storage of data outside HUNT databank, and we have restrictions for researchers for handling of HUNT data files. We have precise information on all data exported to different projects and there are no restrictions regarding data export given approval of applications to HUNT Research Centre.

<http://www.ntnu.edu/hunt/data>

INTERVAL data are available via the European Genotype Phenotype archive (<https://ega-archive.org/datasets/EGAD00001004080>) and other data from the INTERVAL BioResource are available on reasonable request from helpdesk@intervalstudy.org.uk

The data supporting the findings from the kidney tissue investigations are available upon reasonable request to the authors.

Supplemental Note 2: Covariate definition

Covariates included in the regression analyses were: age at the time of examination (years), sex (binary variable), BMI (kg/m²), smoking (current and former/never, self-reported), diabetes (fasting plasma glucose ≥ 126 mg/dl or treatment for diabetes), hypertension (systolic blood pressure ≥ 140 mm Hg or diastolic blood pressure ≥ 90 mm Hg or treatment for hypertension), log-transformed triglycerides (mg/dL), high-density lipoprotein (HDL, mg/dL), and intake of lipid-lowering drugs (yes/no, ATC code C10). There were some minor differences in the variable definitions in the replication studies: in the case of QMDiab, smoking status was a binary

variable based on cotinine detection in blood, diabetes status was assessed by casual glucose (non-fasting) ≥ 200 mg/dl (11.0 mmol/L) or treatment for diabetes or HbA1c $> 6.5\%$, information on lipid-lowering drug intake was not available, and technical covariates (the top 3 principal components from each of the genetic and proteomic datasets) were included; in INTERVAL, self-reported medication history at a later timepoint (two years after proteomic measurement) was used to define hypertension and diabetes, and non-fasting blood samples were used for lipid measurement; in HUNT3, diabetes status was self-reported, lipid-lowering information was not available and non-fasting serum samples were used for lipid measurement.

Supplemental Note 3: Interaction analyses

We conducted interaction analyses in the discovery cohort (KORA) by adding an interaction term to the fully adjusted linear regression model:

$$\ln(\text{eGFR}) \sim \text{protein}_x + \text{interactor} + \text{protein}_x * \text{interactor} + \text{age} + \text{sex} + \text{BMI} + \text{smoking} + \text{diabetes} + \ln(\text{triglycerides}) + \text{HDL} + \text{lipid medication intake} + \text{hypertension}$$

where $\text{protein}_x * \text{interactor}$ is the interaction term, and three interactors were assessed: age, sex, and smoking. One model per protein and per interactor was run. The full set of results for the 80 eGFR-associated proteins identified in KORA F4, the discovery cohort, are provided in Supplementary Table 9. The following interaction terms were significant at a Bonferroni-corrected p-value < 0.05 ($0.05/997$):

- 5 age*protein interactions: b2-Microglobulin, IGFBP-6, FSTL3, JAM-B, Renin

Target	Estimate	95% CI	pval	regression term
b2-Microglobulin	-0.060	-0.07,-0.05	1.12E-30	protein
b2-Microglobulin	-0.003	-0.004,-0.002	1.96E-07	protein:age
IGFBP-6	-0.058	-0.068,-0.048	5.00E-28	protein
IGFBP-6	-0.003	-0.004,-0.001	2.00E-05	protein:age
FSTL3	-0.039	-0.049,-0.029	1.46E-13	protein
FSTL3	-0.002	-0.004,-0.001	3.64E-05	protein:age
JAM-B	-0.035	-0.045,-0.026	1.27E-12	protein
JAM-B	-0.002	-0.004,-0.001	2.19E-05	protein:age
Renin	-0.024	-0.035,-0.013	1.89E-05	protein
Renin	-0.003	-0.004,-0.002	5.61E-07	protein:age

As an example, the significant interaction protein:age term from the model assessing b2-microglobulin suggests the negative association between b2-microglobulin plasma protein levels and log-transformed eGFR is accentuated with age, decreasing by 0.003 [95%CI = -0.004,-

0.002] for each additional year of age ($p = 1.96E-07$) (i.e. change in the slope of standardized protein relatively quantified levels for every one unit increase in age).

In summary, interaction analyses suggested the association between plasma protein levels and log-transformed eGFR varies by age: the association between protein level and log-transformed eGFR for most proteins decreased for each additional year of age. No statistically significant interactions were observed with sex or smoking. Full results for the 80 eGFR-associated proteins identified in KORA F4, the discovery cohort, from all three interaction analyses (age, sex and smoking) are provided in Suppl. Table 9.

Supplemental Note 4: Mendelian randomization analysis

Causality was assessed in the set of 57 proteins whose associations with eGFR showed trans-ethnic replication. Two-sample bidirectional Mendelian randomization¹ (MR) was used to estimate the causal effect of renal function (eGFR as proxy thereof) on plasma protein levels (forward MR) and vice versa (reverse MR, Fig. 1B). All MR analyses were conducted using MRBase as implemented in the R packages *TwoSampleMR* v.4.22 and *MRInstruments* v3.2².

Genetic instruments for eGFR and plasma proteins

We used publicly available genome-wide association (GWA) results for (a) eGFR from the CKDGen consortium, the largest GWA meta-analysis in European-ancestry populations³, and (b) plasma proteins from INTERVAL⁴ and AGES-Reykjavik⁵, the largest studies using SOMAScan with no overlap with the studies included in the CKDGen meta-analysis. All studies had their own research protocols approved by the respective ethics committees and institutional boards, and included written informed consent from the participants.

Analyses in the CKDGen Consortium

A meta-analysis of genome-wide association (GWA) results performed in population-based samples was conducted by the CKDGen Consortium using data from 567,460 participants with European Ancestry (EA) from 85 population-based samples³. Data was imputed to the Haplotype Reference Consortium v1.1 or 1000 Genomes Project. Log transformed serum creatinine-based eGFR was used as the main outcome; some studies also had information on blood urea nitrogen (BUN), an alternative marker of kidney function inversely correlated with eGFR. Renal traits were adjusted for age and sex, and the residuals from these linear regression models were regressed on SNP dosage under an additive genetic model. Study site, genetic principal components, relatedness and other study-specific features were accounted for in the study-specific models. Genome-wide summary statistics were meta-analyzed using a

1
2
3 fixed-effects inverse-variance-weighted approach and made available at [http://ckdgen.imbi.uni-](http://ckdgen.imbi.uni-freiburg.de)
4 [freiburg.de](http://ckdgen.imbi.uni-freiburg.de) ³.
5
6
7

8 *Analyses in the INTERVAL and AGES-Reykjavik studies*

9 In the INTERVAL Study, genotype-protein associations were determined in 3,301 participants
10 with European ancestry⁴. Participation in the INTERVAL study was conditioned on blood
11 donation criteria (exclusion of individuals with history of major disease and/or having recent
12 illnesses). Genotype calling was done using the Affymetrix Axiom UK Biobank array and
13 genotypes imputed using a combined 1000 Genomes Phase 3-UK10K reference panel⁶.
14 Relative protein abundances were measured using an extended version of the SOMAscan
15 assay and natural log-transformed prior to analysis⁴. Log transformed protein levels were
16 adjusted for age, sex, duration between blood draw and processing (binary, ≤ 1 day/ >1 day), and
17 the first three principal components of ancestry using linear regression. Residuals were then
18 rank-inverse normalized and used in linear regression models under an additive genetic model
19 to test association with genetic variants. Summary statistics were made available at
20 <http://www.phpc.cam.ac.uk/ceu/proteins/>.
21
22
23
24
25
26
27

28 Genotype-protein associations were determined in 5,457 participants from the AGES-
29 Reykjavik study with data imputed to the 1000 Genomes Project. Variants within a window of
30 150kbp up- or downstream from the protein coding genes were considered as cis SNPs.
31 Relative concentrations of 4,137 proteins were measured using a custom-designed SOMAscan
32 assay, and a Yeo-Johnson transformation applied prior to analysis. Protein levels were adjusted
33 for age and sex, and linear regression analyses under an additive genetic model were
34 conducted to test for association with genetic variants. Summary statistics for the single
35 strongest SNP (lowest p value) per region from the significant cis-acting SNPs ($p < 5E-08$) were
36 made available⁵.
37
38
39
40
41
42

43 *Instrument selection*

44 A flow diagram presenting a summary of the genetic instrument selection and data
45 harmonization process is shown in Supplementary Fig. 3 and described in detail in
46 Supplementary Table 5.
47

48 For the forward MR, i.e. assessing the effect of renal filtration on protein levels, SNPs
49 showing genome-wide significance with serum creatinine-based eGFR in the CKDGen results
50 were selected as candidate instrumental variables (N = 256). In order to keep genetic signals of
51 renal function and eliminate those more likely reflecting serum creatinine metabolism, only
52 SNPs with a significant association with BUN (one-sided $P < 0.05/256$, N=51) with effect direction
53
54
55
56
57

1
2
3 opposite of that to eGFR (N=211) were kept. 47 SNPs met both criteria and were retained for
4 the next steps. Linkage disequilibrium (LD) clumping was conducted to identify independent
5 variants ($r^2 = 0.01$ and $Kb = 10,000$), resulting in the exclusion of 6 SNPs (rs9828976,
6 rs28817415, rs6484504, rs2472297, rs506000, and rs111827672). Summary-level SNP-protein
7 associations for the remaining 41 eGFR-associated SNPs were extracted from the INTERVAL
8 GWAS results for 47 proteins.
9

10
11
12 For the reverse MR, i.e. interrogating the causal effect of proteins on renal filtration, one
13 to 583 genome-wide significant SNPs for 37 proteins were identified in the INTERVAL GWA
14 results. SNPs were classified as either cis-acting pQTLs (within 1Mb window from the start/end
15 of the protein coding genes) or trans-acting pQTLs (outside the 1Mb window) based on the
16 gene information retrieved with BioMart(Supplementary Table 6). No cis-SNPs were found for 9
17 proteins, whereas one to 581 cis-SNPs were found for 28 proteins. LD clumping was conducted
18 ($r^2 = 0.01$ and $Kb = 10,000$), resulting in a total of one to five independent cis-SNPs per protein.
19 Summary-level SNP-eGFR associations for 28 proteins were then extracted from the CKDGen
20 data (cis-SNP for *UNC5H3* was not found). No SNPs were excluded in the data harmonization
21 process, so that 28 proteins with 1 to 5 SNPs were available for MR.
22
23

24
25
26 The same strategy was followed to identify instruments in the AGES-Reykjavik GWAS .
27
28 Of the 57 replicated eGFR-associated proteins, 20 were not available and eight had no
29 genome-wide significant SNPs. Of the 29 proteins with genome-wide significant SNPs, all had
30 at least one cis-SNP after LD clumping. Instruments for 26 proteins were then extracted from
31 the SNP-outcome data, where SNPs for three proteins/genes (*ST6*, *CD55*, *RETN*) were not
32 found in this dataset. One protein was further excluded in the data harmonization step due to it
33 cis-SNP being a palindromic with intermediate allele frequency (rs28629977 for *EPHA2*).
34
35 Finally, 25 proteins with 1 to 2 SNPs were available for MR.
36
37

38
39
40 In total, 35 unique proteins were tested in this reverse MR direction (i.e. effect of protein
41 on eGFR). MR was performed for 11 proteins using only pGWAS data from INTERVAL and for
42 8 proteins using only pGWAS data from AGES-Reykjavik, whereas 16 proteins were tested
43 using data from both datasets; the overlap between pGWAS datasets is shown in Suppl. Fig. 4.
44 A total of 51 MR analyses were run in the reverse MR direction, thus the threshold for multiple
45 testing correction was set at $p_{\text{Bonferroni}} = 9.43E-04$ ($0.05/51$).
46
47
48
49
50
51

52 *Data harmonization*

53 Genetic effects were aligned to the exposure-increasing allele, and effect alleles, regression
54 coefficients and effect allele frequencies were calculated correspondingly (i.e. effect allele was
55
56
57

1
2
3 coded as other allele, regression coefficient is multiplied by -1 and the new effect allele
4 frequency is obtained after subtracting the old effect allele frequency from one⁷). Formatting and
5 harmonization were done using the *TwoSampleMR* package². Palindromic SNPs with minor
6 allele frequency close to 50% and SNPs with incompatible alleles after data harmonization were
7 excluded⁷. Harmonized datasets used in the MR analyses are available in Supplementary Table
8
9
10
11 8.

12 *Phenotypic variance explained*

13 The proportion of phenotypic variance explained by the SNPs was estimated as

$$14 \beta^2 \left(\frac{2p(1-p)}{\text{var}} \right)$$

15
16
17
18 where β is the SNP effect, p is the effect allele frequency and var the variance of the sex- and
19 age-adjusted phenotype residuals³. For $\log(\text{eGFR})$ it was assumed to be 0.016 on the basis of
20 data from 11,827 European-ancestry participants from the population-based ARIC study⁸.

21 *MR methods*

22 Bidirectional Mendelian randomization was used to investigate the direction of the causal effects
23 between plasma protein levels and renal function (with eGFR as a proxy thereof).

24
25 The primary analysis used inverse variance weighted (IVW) regression for analyses
26 instrumented by two or more SNPs; in this method the coefficient of the gene-outcome
27 association is regressed on the coefficient of the gene-exposure association with the intercept
28 constrained to zero, assuming no directional pleiotropy^{9, 10}. In cases where only one SNP
29 instrumented the analysis, Wald's ratio (coefficient of the gene-outcome association divided by
30 the gene-exposure association) was calculated instead¹⁰.

31
32 To test the three assumptions upon which the validity of MR analyses depend – i.e.
33 relevance (strength of genetic association), independence (specificity of association, i.e. no
34 association with confounders) and exclusion restriction (only associated with outcome through
35 the exposure)¹ – several analyses were conducted. Three further MR methods were used in
36 sensitivity analyses for MR analyses instrumented by more than two SNPs¹¹ (i.e. 47 proteins in
37 the forward MR and 2 proteins in the reverse MR, Supplementary Table 10). MR-Egger
38 regression was used to assess pleiotropy, as this method allows for horizontal pleiotropy and
39 provides an estimate of the unbalanced horizontal pleiotropic effects in its intercept¹². Weighted
40 median¹³ and weighted mode MR¹⁴, methods less sensitive to the presence of invalid
41 instruments and to pleiotropic SNPs behaving as outliers, were also used. Consistency across
42 causal estimates was investigated by performing single SNP analyses. Cochran's Q test was
43 used to test for instrument heterogeneity¹¹ and the MR-PRESSO test was used to assess global
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

1
2
3 heterogeneity¹⁵. MR-PRESSO (*MR-PRESSO* v.2.2 in R) was also used to identify outlier
4 variants (outlier test) and to obtain a causal estimate after exclusion of outliers. Leave-one-out
5 analyses were run to identify individual SNPs potentially biasing the effect estimate, and funnel
6 plots showing the MR estimate against their precision were done to visually check for directional
7 pleiotropy¹¹.
8
9

10
11 Multiple testing was accounted for by the Bonferroni correction, defined as 0.05 divided
12 by the number of total proteins assessed in each MR direction (i.e. 47 in forward MR, 51 in
13 reverse MR). Causal effects were considered robust if these agreed in direction and magnitude
14 across MR methods, were significant at Bonferroni $p < 0.05$ from the IVW estimator and had
15 supporting evidence from MR methods in the sensitivity analyses.
16
17

18 *Instrument specificity*

19 In order to look for further evidence of horizontal pleiotropy, association between our SNPs and
20 other traits were searched for in the GWAS Catalog¹⁶. Eleven of the SNPs instrumenting the
21 forward MR analysis were identified as potentially pleiotropic, given their association with other
22 traits (i.e. blood pressure, lipid dysregulation) as reported in the GWAS Catalog (reported by
23 more than one study at $p < 5E-08$, Supplementary Table 9). Regarding the SNPs used to
24 examine the causal effect of MIA on eGFR, one of the three (rs7937) was reported to be
25 associated with chronic obstructive pulmonary disease (Supplementary Table 9).
26
27
28
29
30
31
32

33 **Supplemental Note 5: Expression analyses in human kidney samples**

34 Gene expression of the protein-coding genes identified in MR (*SPOCK2*, *CA3*, *CST6*, *MIA*) and
35 renal traits was further assessed in (a) data from Nephroseq, a platform of comprehensive renal
36 disease gene expression datasets (N = 458)¹⁷; and (b) RNA-sequencing-derived information on
37 human kidney tissue (N = 427)¹⁸.
38
39
40

41 *Nephroseq analyses*

42 We first used *Nephroseq*¹⁷ (www.nephroseq.org)—a web-based platform for integrative data
43 mining of comprehensive renal disease gene expression datasets—as a resource for
44 association analysis between genes expressed in kidney and eGFR. Gene expression studies
45 on human renal tissue with a minimum of 10 informative individuals were considered for the
46 analyses. A total of 458 kidney samples from four eligible studies by Ju et al.¹⁹ (261 samples),
47 Sampson et al.²⁰ (92 samples), Reich et al.²¹ (31 samples) and Rodwell et al.²² (74 samples)
48 were available for the association analysis. Rodwell et al.²² samples were the only kidney
49 tissues secured from patients without kidney disease. We meta-analysed the measures of
50 association (i.e. Pearson's correlation coefficient) by inverse variance weighted meta-analysis
51
52
53
54
55
56
57

1
2
3 approach using random effect models ²³. The analysis was conducted separately for two
4 different kidney groups. i.e. one from glomerular and cortex samples and the other from
5 tubulointerstitial and medulla samples. Heterogeneity was examined using Cochran's Q test.
6
7

8 *Human kidney tissue collections*

9 We used 427 human kidney samples collected by five studies ¹⁸, namely *moleculAr analysis of*
10 *human kiDney-Manchester renal tlssue pRojEct* (ADMIRE), the TRANScriptome of renal
11 human TissuE Study (TRANSLATE) ²⁴⁻²⁸, and its extension (TRANSLATE-T, 'zero time' pre-
12 implantation biopsy prior to transplantation) ²⁶, *moleculaR analysis of mEchanisms regulating*
13 *gene exPression in post-ischAemic Injury to Renal allograft* (REPAIR) and *Renal gEne*
14 *expreSsion and PredispOsition to cardiovascular and kidNey Disease* (RESPOND) studies.
15
16
17
18

19 In brief, TRANSLATE, ADMIRE and RESPOND studies collected samples from patients with
20 unilateral kidney cancer – the specimen was taken from unaffected by cancer part of the kidney
21 immediately after elective nephrectomy²⁴⁻²⁸. TRANSLATE-T and REPAIR collected pre-
22 implantation kidney biopsies from deceased kidney donors prior to the organ transplantation²⁶.
23 The secured tissue samples from all the studies were immersed immediately in RNAlater or
24 snap-frozen for the purpose of further molecular analysis. All patients were of white-European
25 ethnicity. Further details on this resource of human kidney are described in ¹⁸.
26
27
28
29
30

31 *Histology phenotype data*

32 Histology samples were scored for 6 histological characteristics: glomerular sclerosis,
33 glomerular Bowman's capsule thickening, tubular atrophy, interstitial fibrosis, interstitial
34 inflammation and vascular lesions. All characteristics are scored from 0 to 3, 0 being no
35 /minimal extent of the characteristic present and 3 usually representing extensive damage, as
36 reported before ²⁷.
37
38
39
40

41 *Gene expression data*

42 RNA sequencing data was generated from poly-A selected RNA samples run on Illumina
43 sequencing instruments. Gene expression values were generated by Kallisto in units of
44 transcripts per million (TPM). Before association testing all TPM values are transformed by
45 $\log_2(\text{TPM}+1)$, quantile normalised and standardised using the rank-based inverse normal
46 transformation (expression units are therefore in standard deviations) ²⁶.
47
48
49
50

51 *Association analyses*

52 In all association tests we adjusted for age, sex, BMI, 3 genetic principal components, diabetes
53 and a variable number of surrogate variables (SV) (29 for eGFR and 26 for all histology
54 phenotypes) in line with the computational pipelines developed and reported before ¹⁸. The
55
56
57

1
2
3 surrogate variables allow us to adjust for unmeasured confounding variables, as reported before
4 ¹⁸. All analyses were conducted using R package *limma* with standard multivariate linear
5 regression model (with gene expression as the response), as reported before ¹⁸.
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57

For Peer Review

Legends to Supplementary Tables

Supplementary Table 1. Results from cross-sectional analysis of eGFR and CKD in KORA F4

Results from discovery phase for eGFR and CKD. Provided in the Table are several protein identifiers (TargetFullName, Target, UniProt, EntrezGeneSymbol). Coefficient: linear regression coefficient, SE: standard error, 95% CI lower: lower boundary of 95% confidence interval for the regression coefficient, 95% CI upper: upper boundary of 95% confidence interval for the regression coefficient, N: sample size, Pval: unadjusted p-value, OR: odds ratio, 95% CI OR: 95% confidence interval from odds ratio. The last column shows TRUE if the protein was available in all four studies included (k = 993 in the common set of proteins).

Supplementary Table 2. Validation of proteomic targets from Emilsson et.al., 2018

Provided in the Table are several protein identifiers (TargetFullName, Target, UniProt, EntrezGeneSymbol). Columns starting with "Emilsson_" were extracted from Supplementary Tables 3 and 4 from ⁵; Biological.Matrix: sample used in validation assay.

Supplementary Table 3. Validation of proteomic targets from Sun et.al., 2018

Provided in the Table are several protein identifiers (TargetFullName, Target, UniProt, EntrezGeneSymbol). Columns starting with "Sun_" were extracted from Supplementary Table 3 from ⁴; further experimental details are found in ⁴. The last column shows proteins flagged in terms of cross-reactivity issues: 0 for no binding observed, 1 for binding at least 10-fold weaker than target with product of same gene, 2 for binding at least 10-fold weaker than target with product of different gene, 3 for comparable binding observed to product of same gene and 4 for comparable binding observed to product of different gene.

Supplementary Table 4. Validation of proteomic targets from Williams et.al., 2019

Provided in the Table are several protein identifiers (TargetFullName, Target, UniProt, EntrezGeneSymbol). Columns starting with "Williams_" were extracted from Supplementary Table 3 from ²⁹; further experimental details are found in ²⁹. The last column shows proteins flagged in terms of cross-reactivity issues: 0 for no binding observed, 1 for binding at least 10-fold weaker than target with product of same gene, 2 for binding at least 10-fold weaker than target with product of different gene, 3 for comparable binding observed to product of same gene and 4 for comparable binding observed to product of different gene.

1
2
3 **Supplementary Table 5. Genetic instrument selection and data harmonization (INTERVAL**
4 **and AGES- Reykjavik)**

5
6 Provided in the Table are several protein identifiers (TargetFullName, Target, UniProt,
7 EntrezGeneSymbol). Significant discovery: TRUE if Bonferroni p-value < 0.05 in KORA F4,
8 Significant R1: TRUE if replicated at p-value < 0.05 in meta-analysis and same direction of
9 effect in studies of European Ancestry (i.e. INTERVAL and HUNT), Significant R2: TRUE if
10 replicated at p-value < 0.05 and same direction of effect in non-European study (QMDiab),
11 Significant R3: TRUE if replicated in R1 and R2, Significant eGFR: TRUE if replicated at p-value
12 < 0.05 in trans-ethnic eGFR replication, Significant CKD: TRUE if replicated at p-value < 0.05 in
13 trans-ethnic CKD replication. All columns starting with INTERVAL correspond to information
14 extracted from the proteome GWAS reported in ⁴; PGWAS: whether the protein was reported,
15 IVs in PGWAS: number of genome-wide significant SNPs, N (extracted): number of genome-
16 wide significant SNPs extracted from GWAS data, N (no RSID): SNPs for which no RSID was
17 found, thus dropped from analysis, N (cis) / N (trans): total number of genome-wide significant
18 SNPs of each type, N (LD excl.): number of SNPs excluded in LD pruning, N (cis left) / N (trans
19 left): number of each type of SNPs left after LD pruning, N (left): total number of cis-SNPs
20 qualifying as instrumental variables available in PGWAS data, N (cis in CKDGEN) / N (trans in
21 CKDGEN): number of SNPs found in CKDGen data. The AGES-Reykjavik study was to identify
22 potential instrumental variables for the 57 proteins; reported: found among aforementioned,
23 genome-wide SNPs: TRUE if SNPs with pval < 5E-08 are available, N (cis SNPs) / N (trans):
24 number of genome-wide significant SNPs of each type available, N (cis in CKDGEN) / N (trans
25 in CKDGEN): SNPs found in CKDGEN data, N (cis excl. in harmonization): number of SNPs
26 excluded in data harmonization due to it being a palindromic SNP with intermediate allele
27 frequency, N (cis in CKDGEN final): number of cis-SNPs used in MR, study: last name of study
28 providing GWAS data for identification of IV.
29
30
31
32
33
34
35
36
37
38
39
40
41
42

43 **Supplementary Table 6. Harmonized summary statistics used in MR**

44 Harmonized dataset with summary statistics from ³⁻⁵, as given in the output from the
45 *TwoSampleMR* v.4.22 and *MRInstruments* v3.2² R packages.
46
47
48

49 **Supplementary Table 7. Gene-trait information retrieved from GWAS Catalog for**
50 **pleiotropic SNPs**

51 SNP-trait associations obtained from the GWAS Catalog for all 57 proteins, filtered for SNPs
52 meeting genome-wide significance and SNP-trait associations where the trait may act as a
53 confounder between our examined exposures and outcomes (blood pressure, lipids,
54
55
56
57

1
2
3 cardiovascular disease, serum urate, platelet count, caffeine metabolism) reported by more than
4 one study.
5

6
7 **Supplementary Table 8. Sensitivity analyses with CysC-based eGFR, no adjustment for**
8 **BMI and no adjustment for T2D**
9

10 Results from discovery phase for cystatin C-based eGFR. Provided in the Table are several
11 protein identifiers (TargetFullName, Target). Coefficient: linear regression coefficient, SE:
12 standard error, 95% CI: 95% confidence interval, N: sample size, Pval: unadjusted p-value,
13 Significant eGFR-CysC: TRUE if significant at p-value Bonferroni < 0.05 in this analysis,
14 Significant eGFR-crea: TRUE if significant at p-value Bonferroni < 0.05 in main analysis,
15 Significant no BMI adj.: TRUE if significant at p-value Bonferroni < 0.05 in this analysis,
16 Significant BMI adj.: TRUE if significant at p-value Bonferroni < 0.05 in main analysis,
17 Significant no T2D adj.: TRUE if significant at p-value Bonferroni < 0.05 in this analysis,
18 Significant T2D adj.: TRUE if significant at p-value Bonferroni < 0.05 in main analysis.
19
20
21
22
23
24

25 **Supplementary Table 9. Sensitivity analyses (interaction with age, sex and smoking) in**
26 **KORA F4**
27

28 Results from discovery phase for additional renal phenotypes. Provided in the Table are several
29 protein identifiers (TargetFullName, Target). Analysis: interactor examined (either age, sex and
30 smoking), Estimate: linear regression coefficient, 95%CI: 95% CI of the regression estimate,
31 pval: unadjusted p-value, regression term: regression term from model (i.e. protein:age
32 corresponds to interaction term with age). Significant p-values are highlighted in red.
33
34
35
36

37 **Supplementary Table 10. Results from observational analysis of supplementary renal**
38 **phenotypes (eGFR decline, log(uACR) and MA) in KORA F4**
39

40 Results from discovery phase for additional renal phenotypes. Provided in the Table are several
41 protein identifiers (TargetFullName, Target, UniProt, EntrezGeneSymbol). Coefficient: linear
42 regression coefficient, SE: standard error, 95% CI: 95% confidence interval, N: sample size,
43 Pval: unadjusted p-value, Pval (Bonf.): multiple testing adjusted p-value after Bonferroni. Shown
44 in bold are significant proteins identified in this analysis.
45
46
47
48

49 **Supplementary Table 11. Replication of cross-sectional eGFR-protein associations**

50 Results from trans-ethnic replication phase for the 76 eGFR-associated proteins identified in the
51 discovery sample that were available in the common set of 993 proteins. Provided in the Table
52 are several protein identifiers (TargetFullName, Target), as well as the linear regression
53 estimates from the four studies. Coefficient: linear regression coefficient, SE: standard error,
54
55
56
57

95% CI: 95% confidence interval, N: sample size, Pval: unadjusted p-value, Stouffer's zval: estimated effect from meta-analysis, Stouffer's pval: p-value from meta-analysis, Significant R1: TRUE if meta-analysis of studies of European ancestry (INTERVAL and HUNT) had p-value < 0.05 and the same direction of effect as the discovery, Significant R2: TRUE if replicated at p-value < 0.05 and same direction of effect as the discovery, Significant R3: TRUE if replicated in R1 and R2, Significant CKD: TRUE if replicated at p-value < 0.05 in trans-ethnic CKD replication, Previously associated with kidney function (PMID): PMIDs of literature on the protein and renal function, Reported in podocyte-exosome enriched urine (PMID: 23376485): TRUE if reported in publication, Comparison with Ngo 2020 (PMID: 32958645): "replicated" if one of 43 proteins also described in Ngo, et.al., "not reported" if one of 14 proteins discovered in our study not described in Ngo and colleagues' work, "not checked" if replication not verified. Shown in bold are those proteins showing robust trans-ethnic associations in this analysis.

Supplementary Table 12. Replication of cross-sectional CKD-protein associations

Results from trans-ethnic replication phase for the 34 CKD-associated proteins identified in the discovery sample that were available in the common set of 993 proteins. Provided in the Table are several protein identifiers (TargetFullName, Target), as well as the linear regression estimates from the four studies. Coefficient: linear regression coefficient, SE: standard error, 95% CI: 95% confidence interval, N: sample size, Pval: unadjusted p-value, Significant R1: TRUE if association replicated in HUNT, where replication is defined as p-value<0.05 and same direction of effect as discovery, Significant R2: TRUE if association replicated in QMDiab, where replication is defined as p-value<0.05 and same direction of effect as discovery, Significant R3: TRUE if replicated in R1 and R2, Significant eGFR: TRUE if replicated in trans-ethnic replication. Shown in bold are significant proteins identified in this analysis.

Supplementary Table 13. Extended annotation file (DAVID)

Extended annotation on Gene Ontology terms (biological process – BP_DIRECT, cellular component – CC_DIRECT, molecular function - MF_DIRECT), pathways (BioCarta, KEGG, BBID), functional category (UP_KEYWORDS, UP_SEQ_FEATURE), protein domain/family (SMART, PIR_SUPERFAMILY), disease association (OMIM) as obtained using DAVID v6.8³⁰.

Supplementary Table 14. Gene information from proteins included in MR

Provided in the Table are several protein identifiers (TargetFullName, Target, UniProt, EntrezGeneSymbol), as well as chromosome, gene start/end position, strand and starting/ending positions, annotated to GRCh37 and obtained using BioMartR.

Supplementary Table 15. Phenotypic variance explained by MR instruments

Exposure: log(eGFR) for forward MR, UniProt for reverse MR, followed by several protein identifiers. IVs: SNPs used as instruments, N IVs: number of SNPs instrumenting MR, exposure variance: variance of sex- and age-adjusted phenotype residuals estimated in KORA for plasma proteins, phenotypic variance explained by SNPs: estimated based on the SNP-phenotype effect, effect allele frequency and variance of the sex- and age-adjusted phenotype residual, N: sample size (analyses done in KORA F4).

Supplementary Table 16. Results from MR (IVW, MBE, weighted median and MR-Egger)

id.exposure: ln_eGFR in forward MR (i.e. effect of renal filtration on plasma protein levels) and EntrezGeneSymbol in reverse MR, id.outcome: EntrezGeneSymbol in forward MR (i.e. effect of renal filtration on plasma protein levels) and ln_eGFR in reverse MR, outcome: protein name in forward MR and log(eGFR) in reverse MR, exposure: log(eGFR) in forward MR and protein name in reverse MR, method: MR method used, nsnp: number of SNPs instrumenting the analysis, b: causal estimate, SE: standard error, lower CI: lower bound of 95% CI, upper CI: upper bound of 95% CI, pval: p-value, study: last name and year of publication from studies used for selection of instrumental variables, direction: forward when estimating effect of eGFR on proteins and reverse when estimating effect of proteins on eGFR.

Supplementary Table 17. Sensitivity analyses: Heterogeneity (Cochran's Q test)

id.exposure: ln_eGFR in forward MR (i.e. effect of renal filtration on plasma protein levels) and EntrezGeneSymbol in reverse MR, id.outcome: EntrezGeneSymbol in forward MR (i.e. effect of renal filtration on plasma protein levels) and ln_eGFR in reverse MR, outcome: protein name, exposure: log(eGFR) in forward MR and protein name in reverse MR, method: MR methods used, Q: Cochran's Q, Q_df: degrees of freedom, Q pval: heterogeneity p-value, direction: forward when estimating effect of eGFR on proteins and reverse when estimating effect of proteins on eGFR.

Supplementary Table 18. Sensitivity analyses: Pleiotropy in MR-Egger

id.exposure: ln_eGFR in forward MR (i.e. effect of renal filtration on plasma protein levels) and EntrezGeneSymbol in reverse MR, id.outcome: EntrezGeneSymbol in forward MR (i.e. effect of renal filtration on plasma protein levels) and ln_eGFR in reverse MR, outcome: protein name name in forward MR and log(eGFR) in reverse MR, exposure: log(eGFR) in forward MR and protein name in reverse MR, Egger intercept: intercept term in Egger regression, SE: standard error, Pval: p-value, MR direction: forward when estimating effect of eGFR on proteins and reverse when estimating effect of proteins on eGFR.

Supplementary Table 19. Sensitivity analyses: Leave-one-out analyses

id.exposure: ln_eGFR in forward MR (i.e. effect of renal filtration on plasma protein levels) and EntrezGeneSymbol in reverse MR, id.outcome: EntrezGeneSymbol in forward MR (i.e. effect of renal filtration on plasma protein levels) and ln_eGFR in reverse MR, outcome: protein name, exposure: log(eGFR) in forward MR and protein name in reverse MR, SNP: identifier of specific variant excluded in MR, N: sample size, b: causal estimate, SE: standard error, pval: p-value, direction: forward when estimating effect of eGFR on proteins and reverse when estimating effect of proteins on eGFR.

Supplementary Table 20. Sensitivity analyses: Results from restrictive MR

Comparison of MR results obtained in the main forward analysis (instrumented by 40 SNPs) and sensitivity restrictive MR (after exclusion of eleven SNPs) for multiple methods; id.outcome: EntrezGeneSymbol corresponding to protein coding gene, b: causal estimate, se: standard error, pval: p-value. Highlighted in green are results meeting statistical significance after Bonferroni correction (0.05/47), and highlighted in yellow results meeting nominal significance (pval < 0.05).

Supplementary Table 21. Results from correlation analyses between gene expression and eGFR from Nephroseq datasets

Using datasets curated by *Nephroseq*¹⁷ (www.nephroseq.org), the analysis assessing renal gene expression of the four proteins identified in MR and eGFR was conducted separately for two different kidney groups: one from glomerular and cortex samples and the other from tubulointerstitial and medulla samples, shown on the boxes on the left (dataset_glom_cortex and dataset_tubint_medulla). A total of 458 kidney samples from four eligible studies by Ju et al.¹⁹ (261 samples), Sampson et al.²⁰ (92 samples), Reich et al.²¹ (31 samples) and Rodwell et al.²² (74 samples) were available for the association analysis. The tables on the right show gene_symbol: Entrez Gene symbol, meta_beta: meta-analytic measure of association (i.e. Pearson's correlation coefficient) obtained by inverse variance weighted meta-analysis approach using random effect models²³, mean_se: standard error of estimate, meta_pval: pvalue of meta-analytic estimate, meta_het: Cochran's Q test p value. Highlighted in bold are gene/proteins where a significant (p<0.05) effect was observed.

Supplementary Table 22. Clinical characteristics of studies included in gene expression analyses

Basic clinical characteristics from studies included in the human kidney resource (up to N = 427 kidney samples)¹⁸.

1
2
3 **Supplementary Table 23. Results from multivariate regression analyses on gene**
4 **expression, eGFR and histological characteristic scoring from human kidney resource**

5 Regression analyses on the association between renal traits (eGFR, histologic scores) and
6 gene expression of *SPOCK2*, *MIA*, *CST6* and *CA3*. Regression models included adjustment for
7 age, sex, BMI, 3 genetic principal components, diabetes and a variable number of surrogate
8 variables (29 for eGFR and 26 for all histology phenotypes). Beta: regression estimate, SE:
9 standard error, 95% CI lower and upper: lower and upper bounds of the 95% CI of the
10 regression estimate, P-value: pvalue from regression, N: sample size.
11
12
13
14
15

16 **Supplementary Table 24. Description and biological roles of selected proteins**

17 Protein, description, MW: molecular weight, glomerular filtration / detection in urine: whether the
18 protein is filtrated at the glomeruli and if it has ever been reported in urine, biological role,
19 relevance to kidney function. References used in this table appear at the end of this document.
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

Legends to Supplementary Figures

Supplementary Figure 1. Correlation of serum creatinine variables in KORA F4

Panel A) Quadrant I (top right) shows correlation between NMR-based serum creatinine variable (same as available in INTERVAL) and the standardized Jaffe reaction serum creatinine variable. Quadrants II (top left) and IV (bottom right) show density plot of standardized Jaffe reaction serum creatinine variable and NMR-based variable. Quadrant III (bottom left) shows a scatterplot of both variables. Panel B) Quadrant I (top right) shows correlation between NMR-based eGFR (same as available in INTERVAL) and the standard eGFR variable. Quadrants II (top left) and IV (bottom right) show density plot of both variables. Quadrant III (bottom left) shows a scatterplot of both variables.

Supplementary Figure 2. Correlation between aptamer-based and other measurements for proteins in KORA F4

Scatterplots showing the distribution of the aptamer-based and proximity extension assay measurements of 12 plasma proteins in a subset of the population-based sample studied in the discovery step (KORA F4, N = 174). The correlation between both measurements (Pearson's correlation) and its statistical significance are shown on each plot. Log-transformed normalized data in relative fluorescence units (RFU) from the aptamer-based platform is shown on the x-axis, and log-transformed normalized data in normalized protein expression (NPX) is shown on the y-axis.

Supplementary Figure 3. Genetic instrument selection and data harmonization

Flow diagram showing genetic instrument selection and data harmonization process. Five general steps were taken (identification of genome-wide significant SNPs, additional filtering, LD clumping, extraction of SNP-outcome data, and data harmonization). Number of SNPs identified, excluded and kept for analysis are provided for each step are given.

Supplementary Figure 4. Protein overlap in pGWAS datasets used in reverse direction of MR

Venn diagram showing the sets of proteins identified in the pGWAS datasets from INTERVAL and AGES-Reykjavik studies. The left circle shows the set of proteins for which the MR analysis was performed using INTERVAL pGWAS summary statistics, whereas the right circle shows the set of proteins for which the MR analysis was performed using AGES-Reykjavik summary statistics. In total, 35 unique proteins were tested in the reverse direction of the MR analyses (i.e. effect of protein on eGFR). MR was performed for 11 proteins using only pGWAS data from

INTERVAL and for 8 proteins using only pGWAS data from AGES-Reykjavik, whereas 16 proteins were tested using data from both datasets.

Supplementary Figure 5. Proteins and log(eGFR) distribution in discovery dataset

Scatterplots showing log-transformed estimated glomerular filtration rate (eGFR) on the x-axis and protein plasma concentrations in relative-fluorescence units (RFU) on the y-axis. Pearson's correlation coefficients and their corresponding p-values are shown in the plot. Shown are the top 10 proteins. Protein-log(eGFR) observations are color coded in agreement with chronic kidney disease (CKD) staging³¹: G1 as normal or high GFR (GFR > 90 mL/min) in blue, G2 as mild CKD (GFR = 60-89 mL/min) in green, and G3 as moderate CKD (GFR = 30-59 mL/min) in red. ARMEL: Cerebral dopamine neurotrophic factor, TNF SR-I: Tumor necrosis factor receptor superfamily member 1A, DAN: Neuroblastoma suppressor of tumorigenicity 1, RGMB: RGM domain family member B, FSTL3: Follistatin-related protein 3, JAM-B: Junctional adhesion molecule B.

Supplementary Figure 6. Cross sectional results for eGFR-protein associations across studies

Volcano plots from all included studies: panel A) data from N = 995 individuals from KORA F4, panel B) data from N = 930 individuals from HUNT3, panel C) data from N = 623 individuals from INTERVAL, and panel D) data from N = 334 individuals from QMDiab. The linear regression coefficient for the protein term is shown on the x-axis, and statistical significance as $-\log(p\text{value})$ on the y-axis. Labeled and colored in red are significant associations at $p < 0.05$ in each individual study.

Supplementary Figure 7. Proteins and log(eGFR) distribution in discovery dataset after CKD exclusion

Scatterplots showing log-transformed estimated glomerular filtration rate (eGFR) on the x-axis and protein plasma concentrations in relative-fluorescence units (RFU) on the y-axis in the discovery dataset (KORA F4) after exclusion of 38 individuals with CKD. Pearson's correlation coefficients and their corresponding p-values are shown in the plot. Shown are the top 10 proteins. Protein-log(eGFR) observations are color coded in agreement with chronic kidney disease (CKD) staging³¹: G1 as normal or high GFR (GFR > 90 mL/min) in blue, G2 as mild CKD (GFR = 60-89 mL/min) in green, and G3 as moderate CKD (GFR = 30-59 mL/min) in red. ARMEL: Cerebral dopamine neurotrophic factor, TNF SR-I: Tumor necrosis factor receptor superfamily member 1A, DAN: Neuroblastoma suppressor of tumorigenicity 1, RGMB: RGM

1
2
3 domain family member B, FSTL3: Follistatin-related protein 3, JAM-B: Junctional adhesion
4 molecule B.
5

6
7 **Supplementary Figure 8. Correlation between Z-values for eGFR-protein associations**
8 **across studies**
9

10 Correlation between z-values (regression coefficient divided by its standard error) across all
11 included studies was calculated for the set of proteins significantly associated with eGFR in the
12 discovery study (k=76 proteins at $p < \text{Bonferroni}$ in discovery). The plot shows pairwise
13 scatterplots and correlations between Z-values between each pair of studies, as well as four
14 plots of the density of the z-value distribution of each study in the diagonal of the plot matrix.
15
16
17

18
19 **Supplementary Figure 9. Tissue expression of 57 eGFR-associated proteins**
20 **(ProteomeDB)**
21

22 Tissue expression of 57 eGFR-associated proteins (ProteomeDB). Heatmap showing
23 expression of proteins as rows and biological sources as columns, respectively. The
24 dendrograms show the results of hierarchical clustering of proteins and biological sources.
25 Protein expression values were produced by the MS1 quantification technique and expression
26 values estimated by the iBAQ approach. Highlighted with yellow is the column corresponding to
27 kidney tissue. The data presented in this figure were generated through a multi-protein query
28 using the UniProt IDs on the ProteomicsDB Analytics Toolbox portal on Oct. 2, 2020 from:
29 <https://www.proteomicsdb.org/proteomicsdb/#analytics/expressionHeatmap>
30
31
32
33
34

35
36 **Supplementary Figure 10. Tissue expression of 56 eGFR-associated protein coding**
37 **genes (ProteomeDB)**
38

39 Heatmap showing expression of 56 proteins-coding genes (the set of 57 proteins identified in as
40 trans-ethnically associated with eGFR are the product of 56 genes) as rows and biological
41 sources as columns, respectively. The dendrograms show the results of hierarchical clustering
42 of proteins and biological sources. Panel A shows data from RNA-seq experiments, and panel B
43 data from microarray experiments. Highlighted with yellow is the column corresponding to
44 kidney tissue. The data presented in this figure were generated through a multi-protein query
45 using the UniProt IDs on the ProteomicsDB Analytics Toolbox portal on Oct. 2, 2020 from:
46 <https://www.proteomicsdb.org/proteomicsdb/#analytics/expressionHeatmap>
47
48
49
50
51
52
53
54
55
56
57

1
2
3 **Supplementary Figure 11. Expression of 56 eGFR-associated protein coding genes**
4 **across tissues (GTEx)**
5

6 Gene expression of 57 eGFR-associated proteins across tissues (GTEx). This heatmap
7 provides a qualitative measure of relative expression across human tissues, shown in
8 Transcripts Per Million (TPM) (<https://doi.org/10.1093/bioinformatics/btp692>). Values across
9 tissues may not be compared due to the differences in sample normalization across the diverse
10 set of tissues represented in GTEx (<https://doi.org/10.1016/j.cell.2012.10.012>). The columns
11 corresponding to cortex and medulla kidney tissue samples are highlighted in orange. The data
12 presented in this plot was generated on Oct. 2, 2020 through a multi-gene query on the GTEx
13 portal <https://www.gtexportal.org/home/multiGeneQueryPage>
14
15
16
17
18

19 **Supplementary Figure 12. Protein-protein interaction network of 57 replicated eGFR-**
20 **associated proteins**
21

22 Protein-protein interaction network of 57 replicated eGFR-associated proteins. Colored proteins
23 are query proteins, no specific color-coding is applied. Known interactions: cyan blue edges are
24 retrieved from curated databases, pink edges are experimentally determined; predicted
25 interactions: bright green edges are retrieved from gene neighborhood, red edges from gene
26 fusions and navy blue from gene co-occurrence; others: golden edges are retrieved from text
27 mining, black from co-expression and light purple from protein homology.
28
29
30
31
32

33 **Supplementary Figures 13-19. Forward MR results for effects of eGFR on proteins.**

34 Panel A) Forest plot showing IVW causal estimates following a leave-one-out approach (ie.
35 presented is the casual estimate using the leave-one-out approach, when the given SNP is not
36 included in the analysis). Panel B) Funnel plot showing the ratio estimate for each variant on the
37 x-axis and its square root precision on the y-axis, where asymmetry suggests directional
38 pleiotropy. Vertical lines represent the causal estimates obtained in each method, color coded
39 as in panel C.
40
41
42
43

44 **Supplementary Figure 20. Reverse MR analysis for MIA-eGFR**

45 Panel A) Forest plot showing individual contributions (x-axis) of each SNP instrumenting the
46 analysis (y-axis), followed by the pooled MR estimates obtained with all MR methods. IVW:
47 inverse variance-weighted MR. Panel B) Scatter plot of summary data estimates for the
48 associations of 40 SNPs with log(eGFR) (x-axis) and plasma proteins (y-axis). The lines
49 correspond to the slopes from the IVW (shown in light blue), weighted median (shown in light
50 green), MR-Egger (shown in dark blue) and weighted mode (shown in dark green). Panel C)
51 Funnel plot showing IVW causal estimates following a leave-one-out approach.
52
53
54
55
56
57

Supplementary Figure 21. SPOCK2 gene expression in renal tissue from 26 CKD patients

Scatterplot showing gene expression from microdissected tubulointerstitial components of human renal biopsies from 26 individuals with CKD at different disease stages (I-IV)¹⁹. eGFR is shown in the x-axis and renal *SPOCK2* gene expression in the y-axis. Stages of chronic kidney disease are color coded following the KDIGO's GFR categories: stage 1: normal or high GFR (≥ 90 ml/min/ 1.73m²); stage II: mildly decreased (60-29 ml/min/ 1.73m²); stage III: mild to moderately decreased (59-30 ml/min/ 1.73m²); stage IV: severely decreased (15-29 ml/min/ 1.73m²); stage V: kidney failure (15 ml/min/ 1.73m²). Shown in blue is the regression line corresponding to the eGFR ~ *SPOCK2* expression model.

For Peer Review

References

1. Pierce, BL, Burgess, S: Efficient design for Mendelian randomization studies: subsample and 2-sample instrumental variable estimators. *American journal of epidemiology*, 178: 1177-1184, 2013.
2. Hemani, G, Zheng, J, Elsworth, B, Wade, KH, Haberland, V, Baird, D, et al.: The MR-Base platform supports systematic causal inference across the human phenome. *eLife*, 7, 2018.
3. Wuttke, M, Li, Y, Li, M, Sieber, KB, Feitosa, MF, Gorski, M, et al.: A catalog of genetic loci associated with kidney function from analyses of a million individuals. *Nature genetics*, 51: 957-972, 2019.
4. Sun, BB, Maranville, JC, Peters, JE, Stacey, D, Staley, JR, Blackshaw, J, et al.: Genomic atlas of the human plasma proteome. *Nature*, 558: 73-79, 2018.
5. Emilsson, V, Ilkov, M, Lamb, JR, Finkel, N, Gudmundsson, EF, Pitts, R, et al.: Co-regulatory networks of human serum proteins link genetics to disease. *Science (New York, NY)*, 361: 769-773, 2018.
6. Astle, WJ, Elding, H, Jiang, T, Allen, D, Ruklisa, D, Mann, AL, et al.: The Allelic Landscape of Human Blood Cell Trait Variation and Links to Common Complex Disease. *Cell*, 167: 1415-1429.e1419, 2016.
7. Hartwig, FP, Davies, NM, Hemani, G, Davey Smith, G: Two-sample Mendelian randomization: avoiding the downsides of a powerful, widely applicable but potentially fallible technique. *International journal of epidemiology*, 45: 1717-1726, 2016.
8. Gorski, M, van der Most, PJ, Teumer, A, Chu, AY, Li, M, Mijatovic, V, et al.: 1000 Genomes-based meta-analysis identifies 10 novel loci for kidney function. *Scientific reports*, 7: 45040, 2017.
9. Burgess, S, Butterworth, A, Thompson, SG: Mendelian randomization analysis with multiple genetic variants using summarized data. *Genetic epidemiology*, 37: 658-665, 2013.
10. Teumer, A: Common Methods for Performing Mendelian Randomization. *Frontiers in cardiovascular medicine*, 5: 51, 2018.
11. Burgess, S, Bowden, J, Fall, T, Ingelsson, E, Thompson, SG: Sensitivity Analyses for Robust Causal Inference from Mendelian Randomization Analyses with Multiple Genetic Variants. *Epidemiology*, 28: 30-42, 2017.

12. Bowden, J, Davey Smith, G, Burgess, S: Mendelian randomization with invalid instruments: effect estimation and bias detection through Egger regression. *International journal of epidemiology*, 44: 512-525, 2015.
13. Bowden, J, Davey Smith, G, Haycock, PC, Burgess, S: Consistent Estimation in Mendelian Randomization with Some Invalid Instruments Using a Weighted Median Estimator. *Genetic epidemiology*, 40: 304-314, 2016.
14. Hartwig, FP, Davey Smith, G, Bowden, J: Robust inference in summary data Mendelian randomization via the zero modal pleiotropy assumption. *International journal of epidemiology*, 46: 1985-1998, 2017.
15. Verbanck, M, Chen, C-Y, Neale, B, Do, R: Detection of widespread horizontal pleiotropy in causal relationships inferred from Mendelian randomization between complex traits and diseases. *Nature genetics*, 50: 693-698, 2018.
16. Buniello, A, MacArthur, JAL, Cerezo, M, Harris, LW, Hayhurst, J, Malangone, C, et al.: The NHGRI-EBI GWAS Catalog of published genome-wide association studies, targeted arrays and summary statistics 2019. *Nucleic acids research*, 47: D1005-d1012, 2019.
17. Martini, S, Eichinger, F, Nair, V, Kretzler, M: Defining human diabetic nephropathy on the molecular level: Integration of transcriptomic profiles with biological knowledge. *Reviews in Endocrine and Metabolic Disorders*, 9: 267-274, 2008.
18. Jiang, X, Eales, JM, Scannali, D, Nazgiewicz, A, Prestes, P, Maier, M, et al.: Hypertension and renin-angiotensin system blockers are not associated with expression of angiotensin-converting enzyme 2 (ACE2) in the kidney. *European Heart Journal*, 2020.
19. Ju, W, Nair, V, Smith, S, Zhu, L, Shedden, K, Song, P, et al.: Tissue transcriptome-driven identification of epidermal growth factor as a chronic kidney disease biomarker. *Science translational medicine*, 7: 316ra193, 2015.
20. Sampson, MG, Robertson, CC, Martini, S, Mariani, LH, Lemley, KV, Gillies, CE, et al.: Integrative Genomics Identifies Novel Associations with APOL1 Risk Genotypes in Black NEPTUNE Subjects. *Journal of the American Society of Nephrology : JASN*, 27: 814-823, 2016.
21. Reich, HN, Tritchler, D, Cattran, DC, Herzenberg, AM, Eichinger, F, Boucherot, A, et al.: A molecular signature of proteinuria in glomerulonephritis. *PLoS one*, 5: e13451, 2010.
22. Rodwell, GE, Sonu, R, Zahn, JM, Lund, J, Wilhelmy, J, Wang, L, et al.: A transcriptional profile of aging in the human kidney. *PLoS biology*, 2: e427, 2004.
23. Balduzzi, S, Rücker, G, Schwarzer, G: How to perform a meta-analysis with R: a practical tutorial. *Evidence Based Mental Health*, 22: 153-160, 2019.

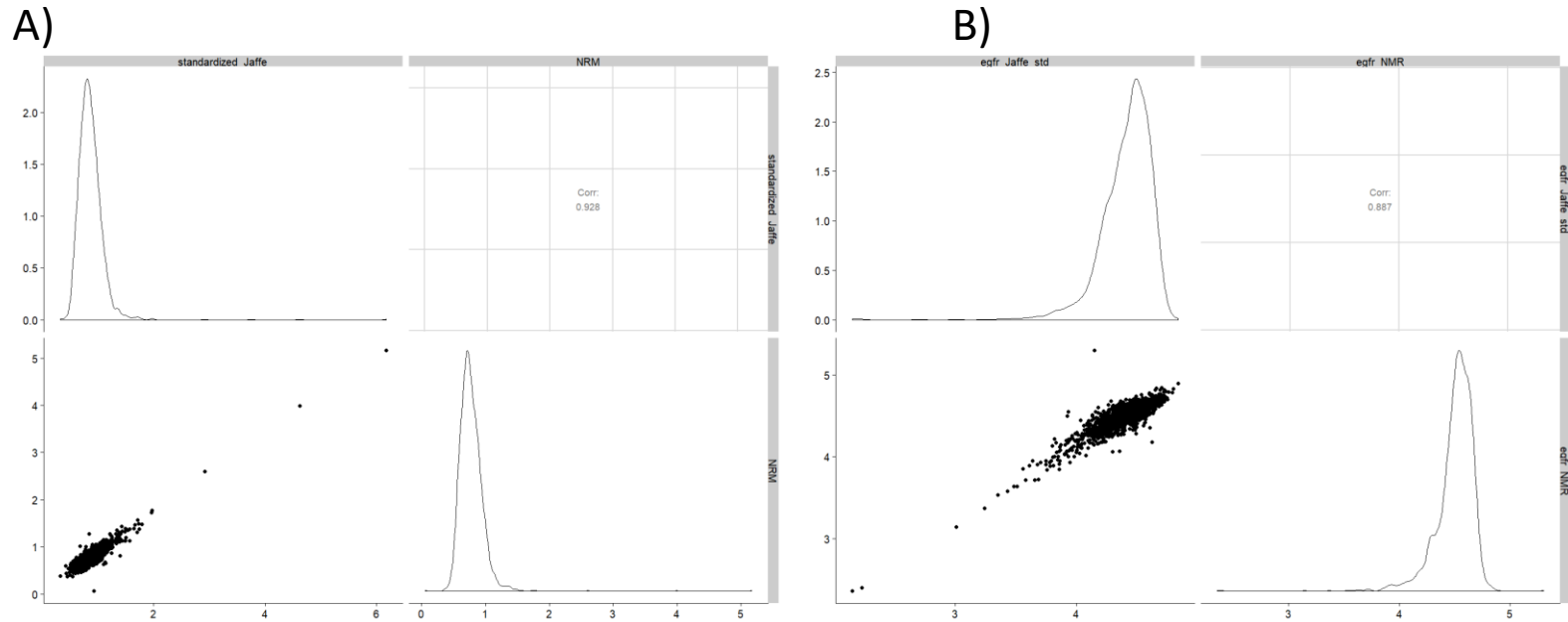
- 1
2
3 24. Marques, FZ, Romaine, SP, Denniff, M, Eales, J, Dormer, J, Garrelds, IM, et al.: Signatures
4 of miR-181a on the Renal Transcriptome and Blood Pressure. *Molecular medicine*
5 (Cambridge, Mass), 21: 739-748, 2015.
6
7
8 25. Tomaszewski, M, Eales, J, Denniff, M, Myers, S, Chew, GS, Nelson, CP, et al.: Renal
9 Mechanisms of Association between Fibroblast Growth Factor 1 and Blood Pressure.
10 *Journal of the American Society of Nephrology : JASN*, 26: 3151-3160, 2015.
11
12 26. Xu, X, Eales, JM, Akbarov, A, Guo, H, Becker, L, Talavera, D, et al.: Molecular insights into
13 genome-wide association studies of chronic kidney disease-defining traits. *Nature*
14 *communications*, 9: 4800, 2018.
15
16 27. Rowland, J, Akbarov, A, Eales, J, Xu, X, Dormer, JP, Guo, H, et al.: Uncovering genetic
17 mechanisms of kidney aging through transcriptomics, genomics, and epigenomics.
18 *Kidney Int*, 95: 624-635, 2019.
19
20 28. Morris, AP, Le, TH, Wu, H, Akbarov, A, van der Most, PJ, Hemani, G, et al.: Trans-ethnic
21 kidney function association study reveals putative causal genes and effects on kidney-
22 specific disease aetiologies. *Nature communications*, 10: 29, 2019.
23
24 29. Williams, SA, Kivimaki, M, Langenberg, C, Hingorani, AD, Casas, JP, Bouchard, C, et al.:
25 Plasma protein patterns as comprehensive indicators of health. *Nature medicine*, 25:
26 1851-1857, 2019.
27
28 30. Huang, DW, Sherman, BT, Lempicki, RA: Systematic and integrative analysis of large gene
29 lists using DAVID bioinformatics resources. *Nature Protocols*, 4: 44-57, 2009.
30
31 31. Hill, NR, Fatoba, ST, Oke, JL, Hirst, JA, O'Callaghan, CA, Lasserson, DS, et al.: Global
32 Prevalence of Chronic Kidney Disease - A Systematic Review and Meta-Analysis. *PLoS*
33 *one*, 11: e0158765, 2016.
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41

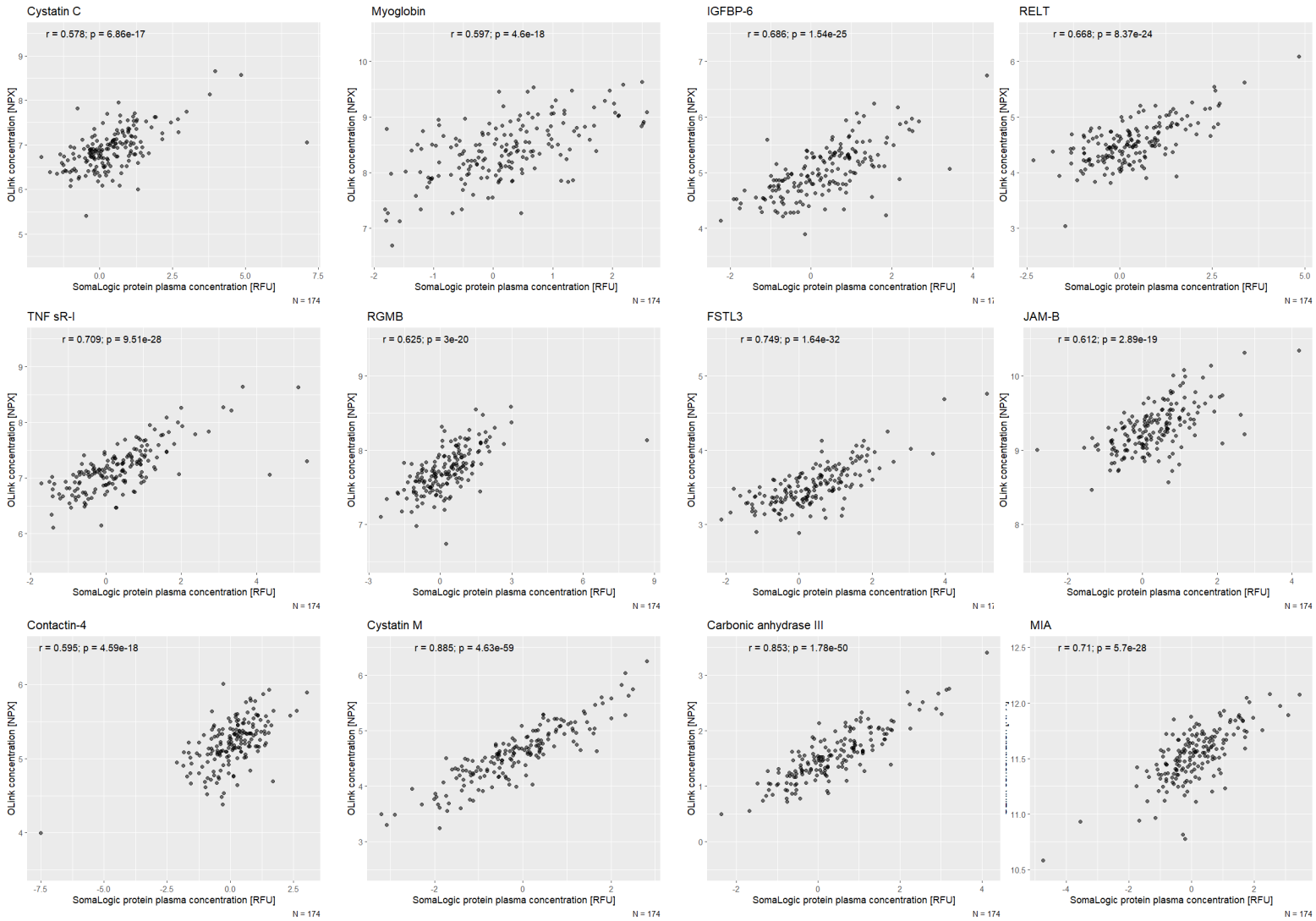
Supplementary Figures

Plasma proteomics of renal function: a trans-ethnic meta-analysis and Mendelian randomization study

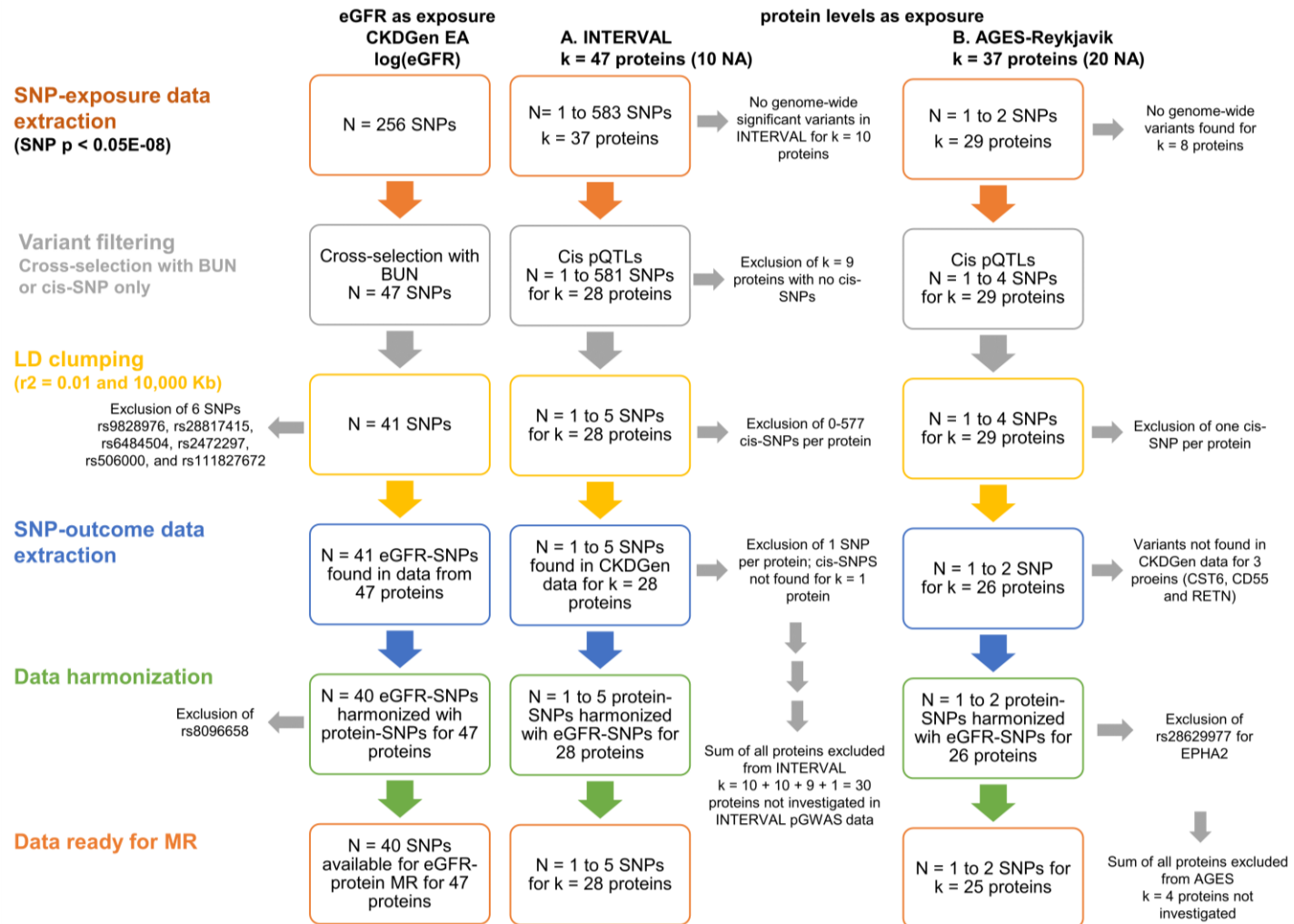
Suppl. Fig. 1. Correlation of serum creatinine variables in KORA F4



Suppl. Fig. 2. Correlation between plasma proteomic measurements in KORA F4



Suppl. Fig. 3. Genetic instrument selection and data harmonization for MR

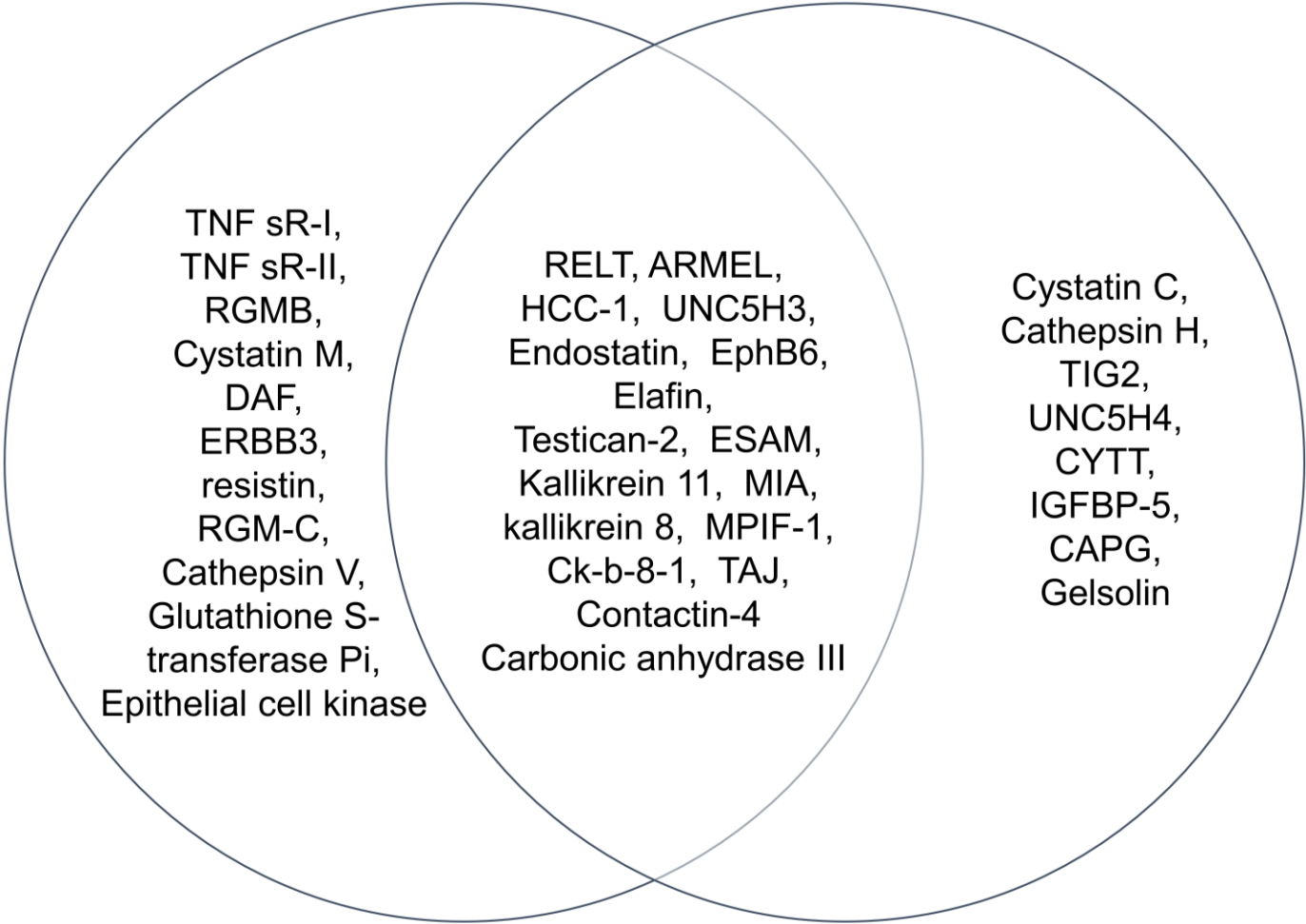


1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41

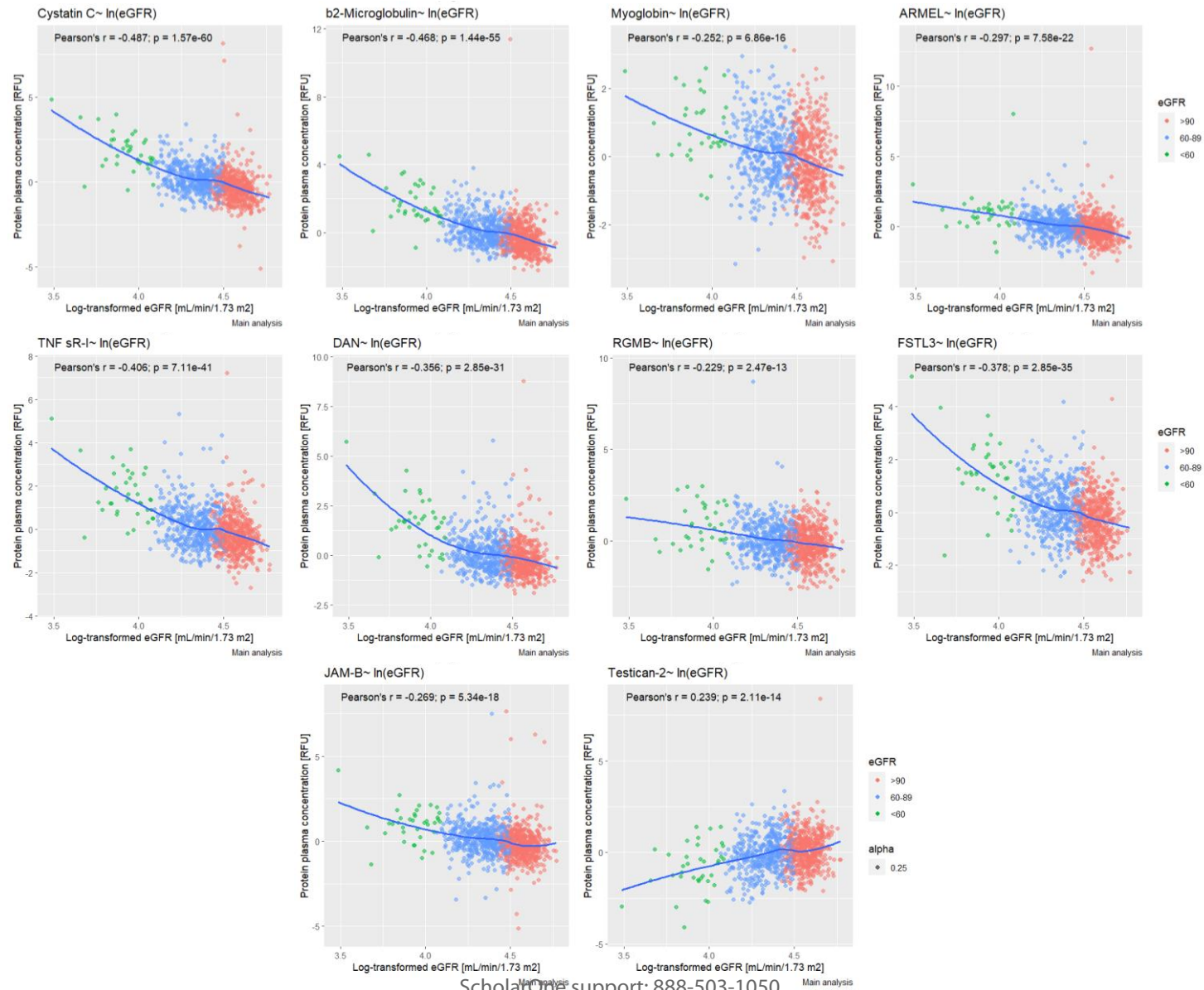
Suppl. Fig. 4. Protein overlap in pGWAS datasets used in reverse direction of MR

Proteins for which the MR analysis was performed using **INTERVAL pGWAS** summary statistics

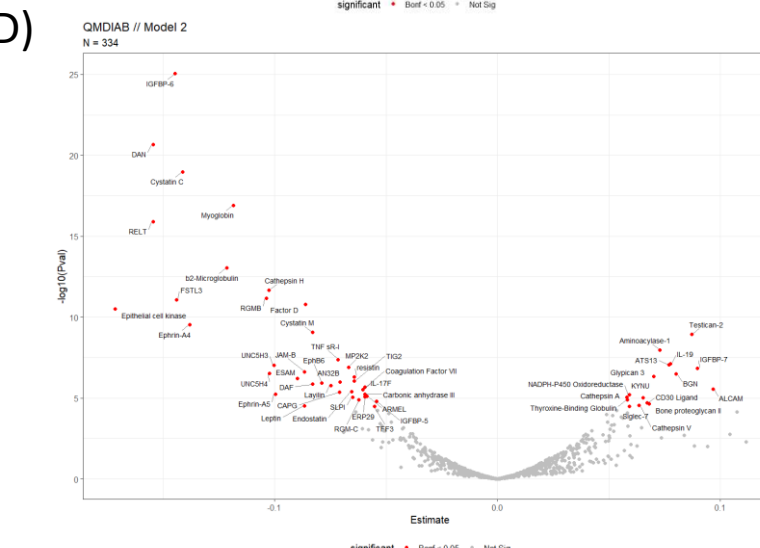
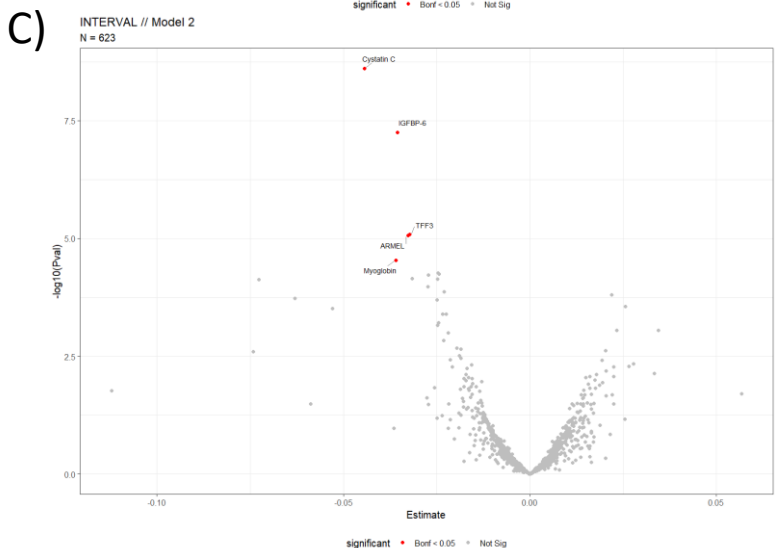
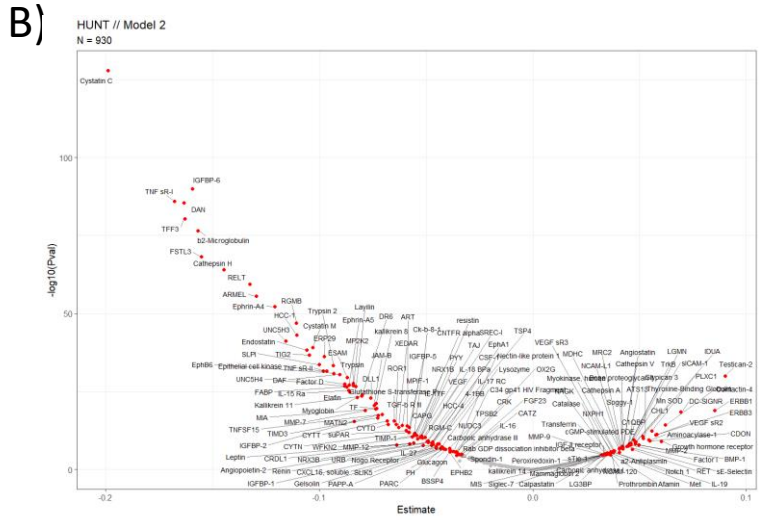
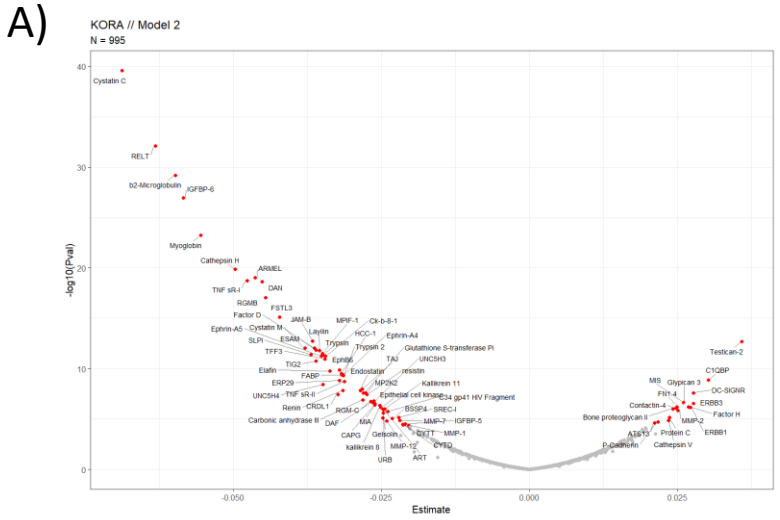
Proteins for which the MR analysis was performed using **AGES-Reykjavik pGWAS** summary statistics



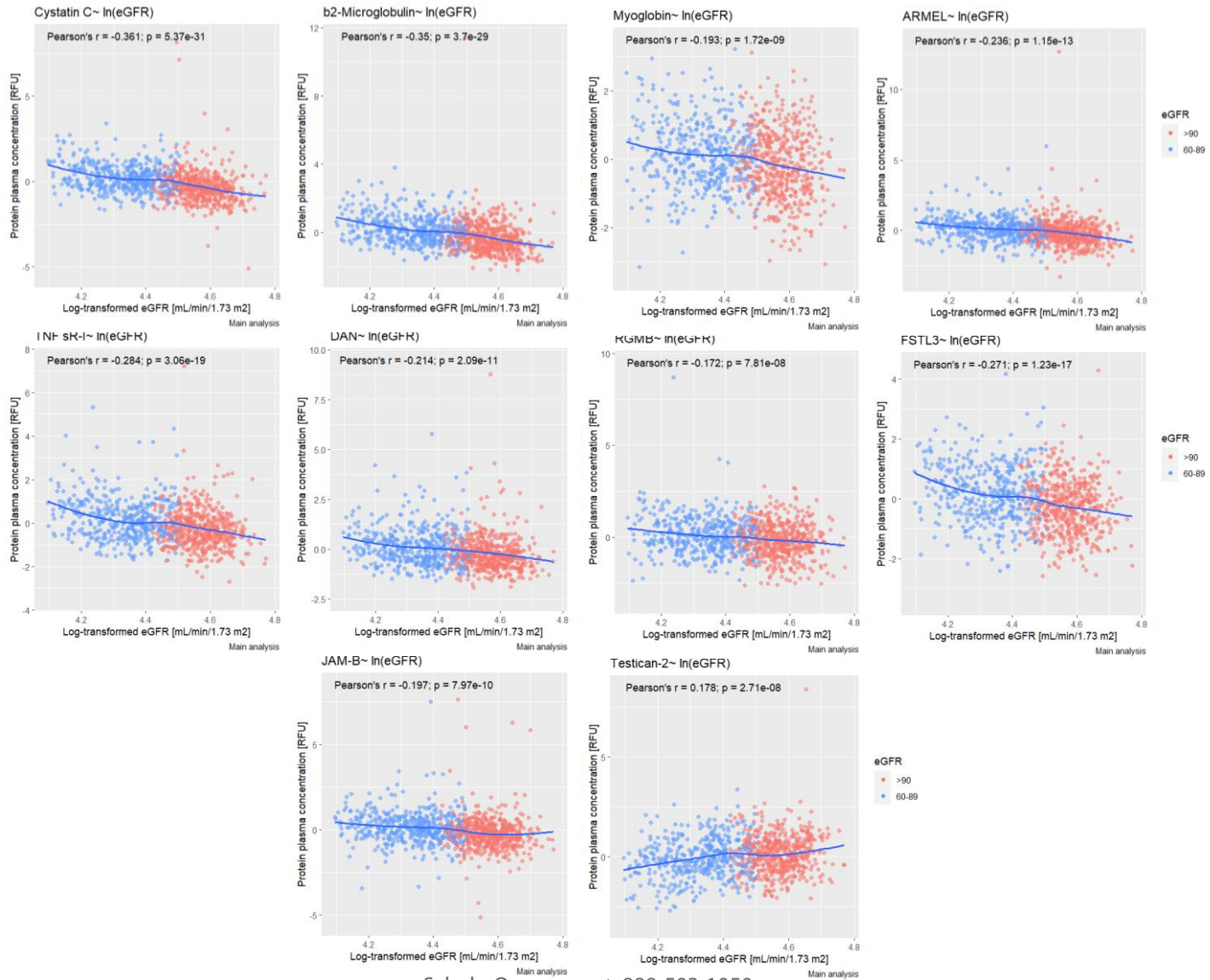
Suppl. Fig. 5. Proteins and log(eGFR) distribution in discovery dataset



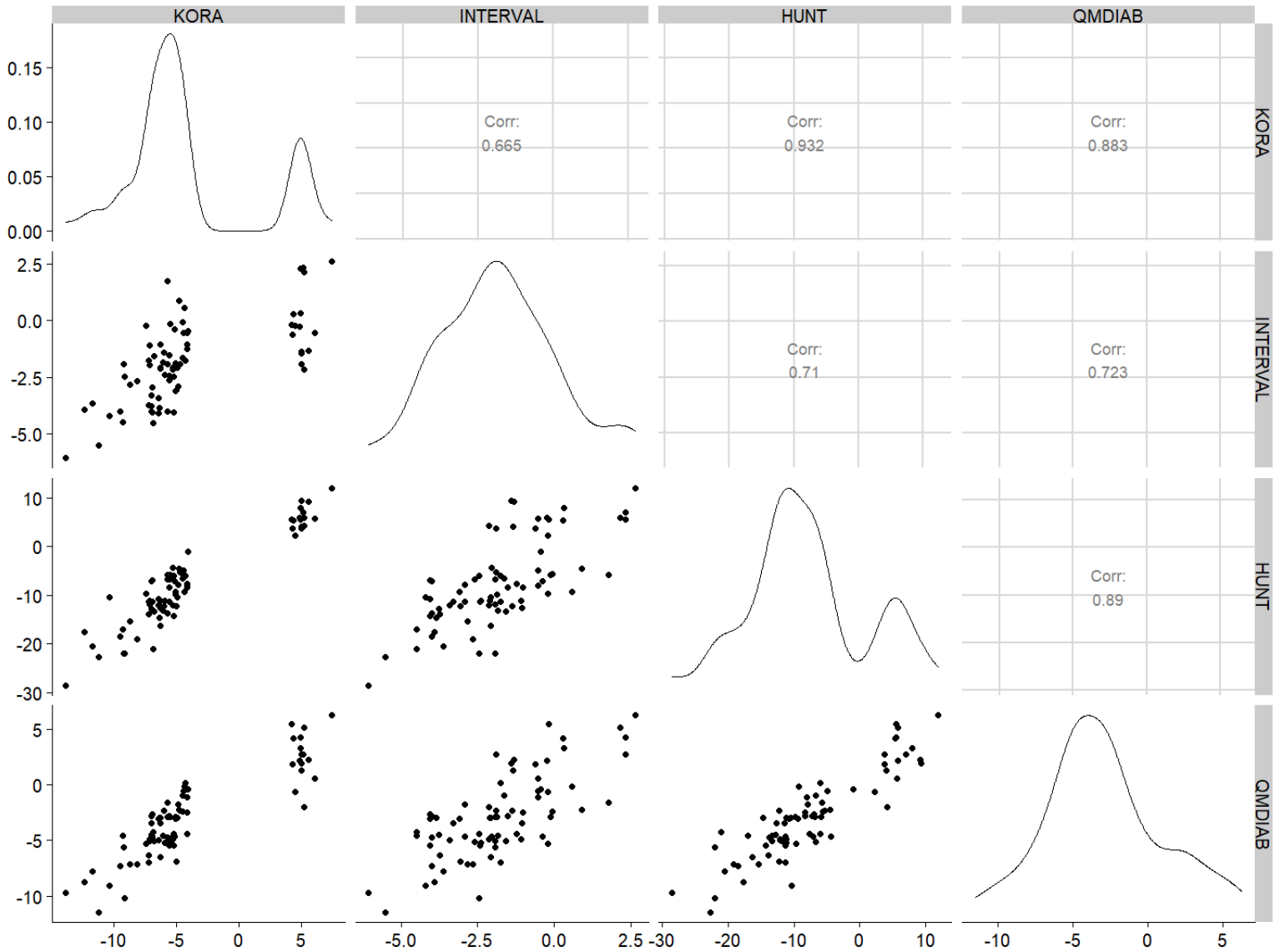
Suppl. Fig. 6. Cross sectional results for eGFR-protein associations across studies



Suppl. Fig. 7. Proteins and log(eGFR) distribution in discovery dataset after CKD exclusion



Suppl. Fig. 8. Correlation between Z-values for eGFR-protein associations across studies



1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41

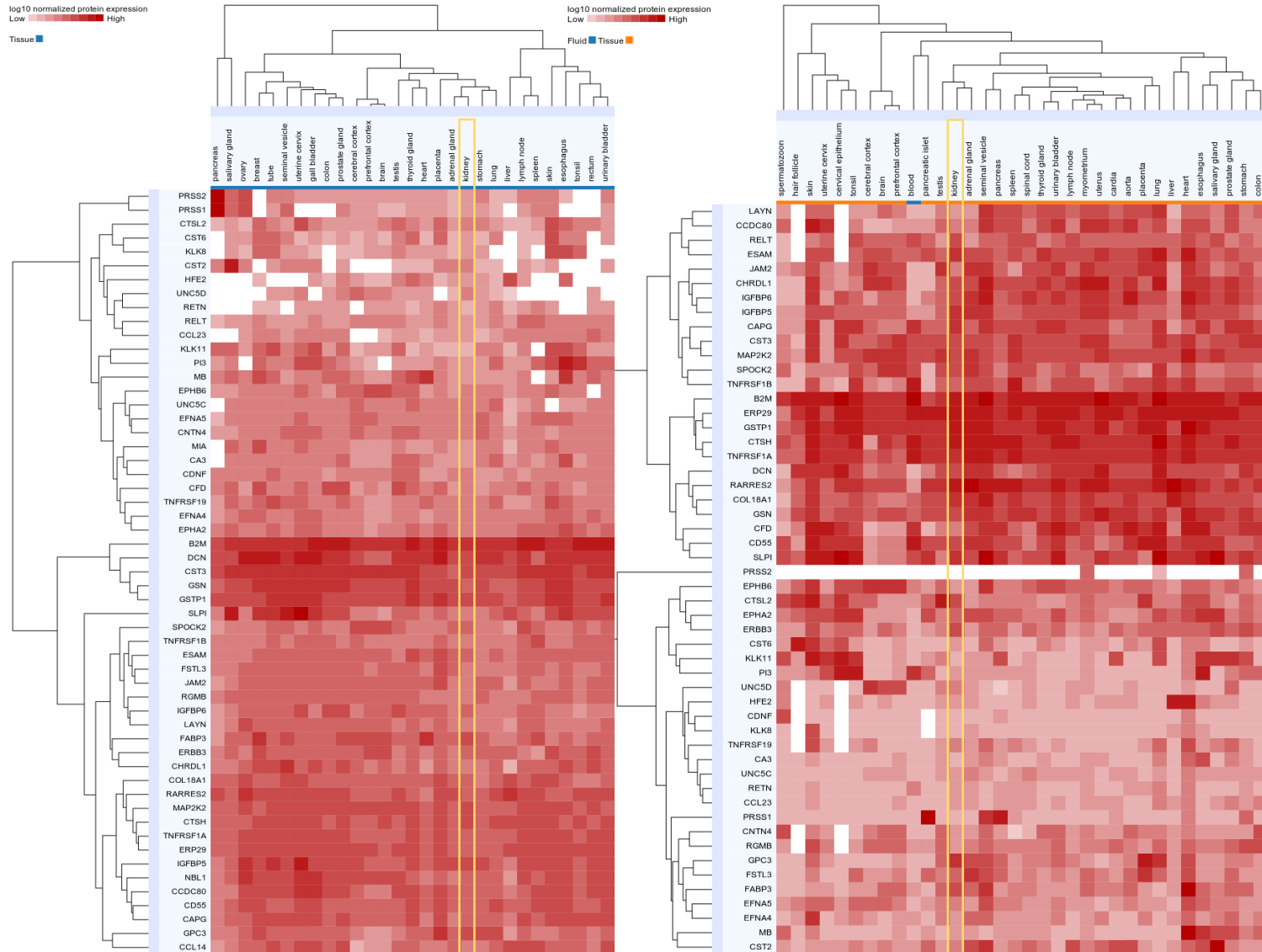
Suppl. Fig. 9. Tissue expression of 56 eGFR-associated proteins (ProteomeDB)

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41



Suppl. Fig. 10. Tissue expression of 56 eGFR-associated protein coding genes (ProteomeDB)

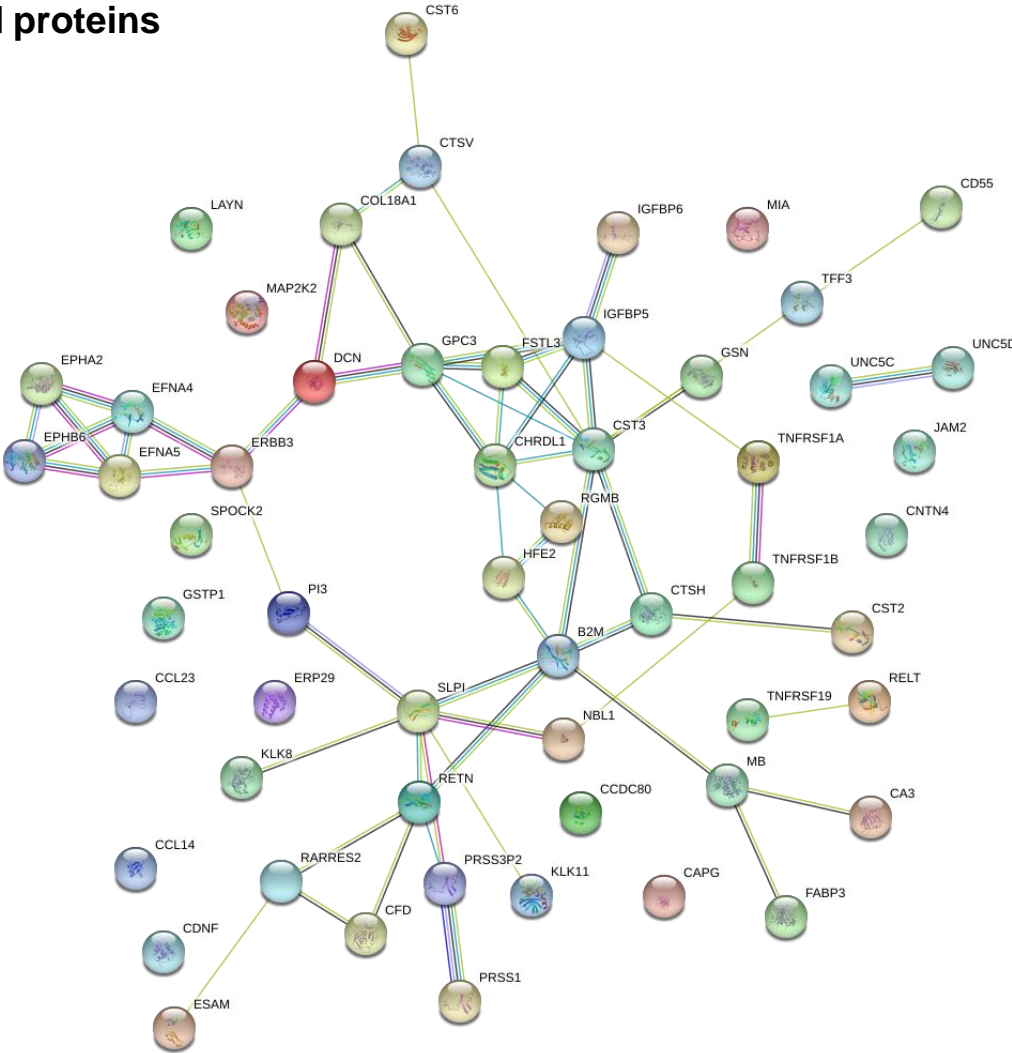
1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41



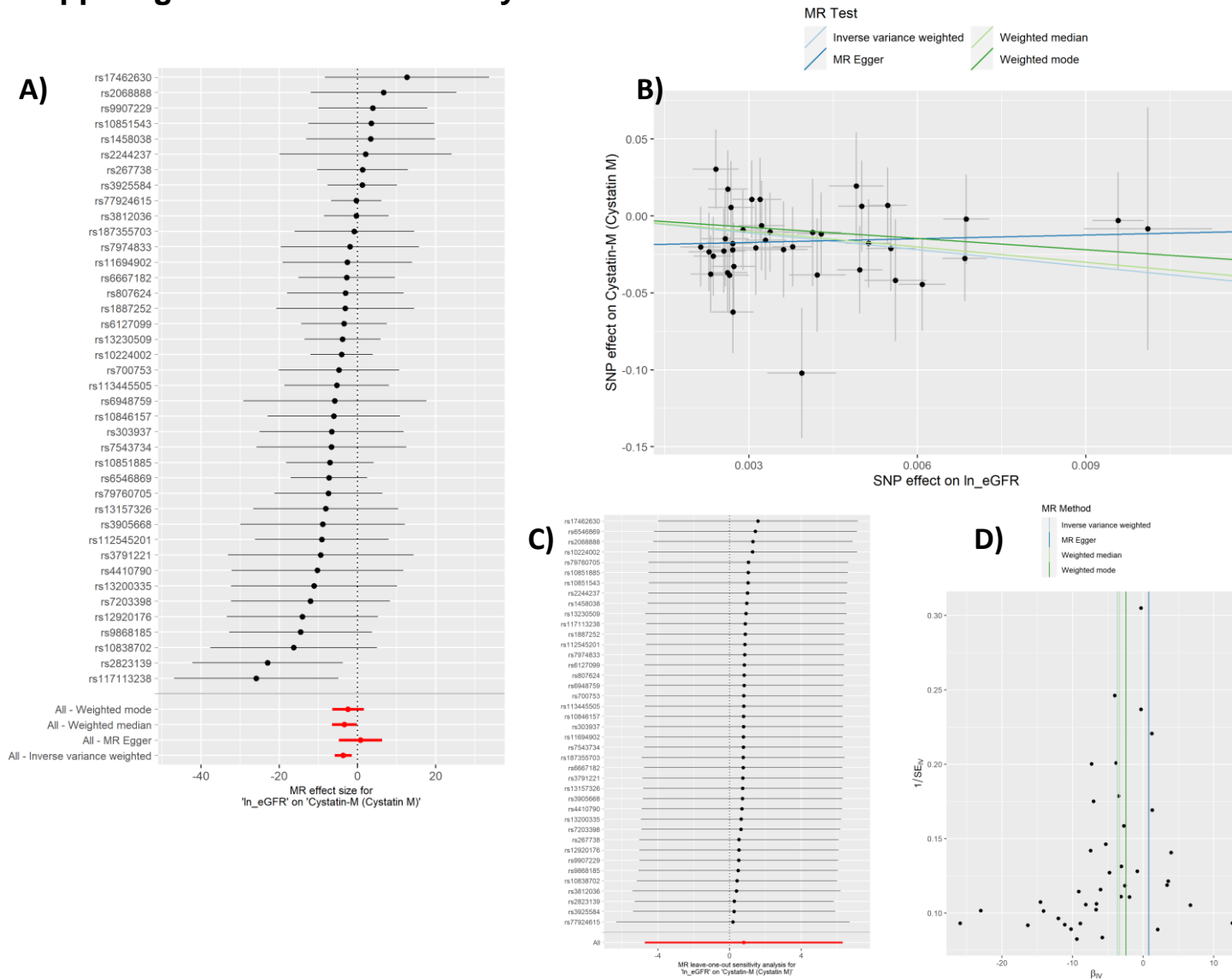
ScholarOne support: 888-503-1050

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41

Suppl. Fig. 12. Protein-protein interaction network of 57 replicated eGFR-associated proteins

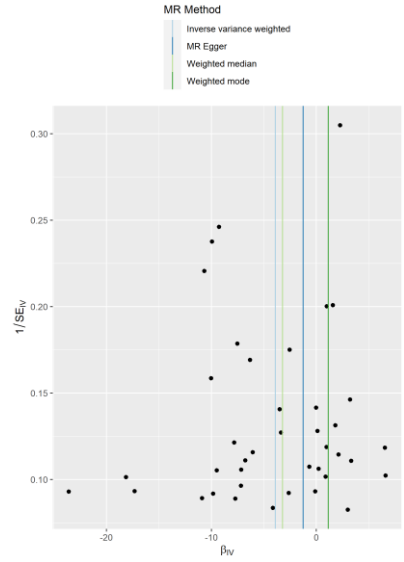
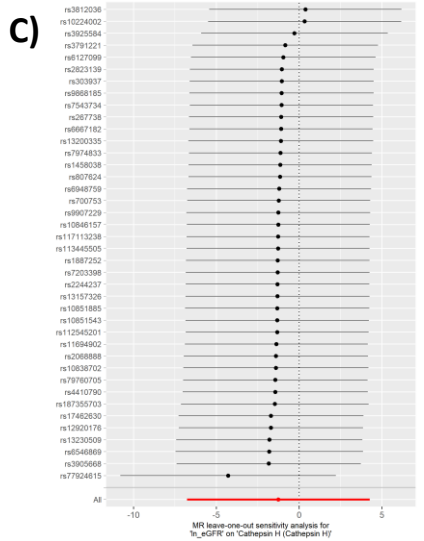
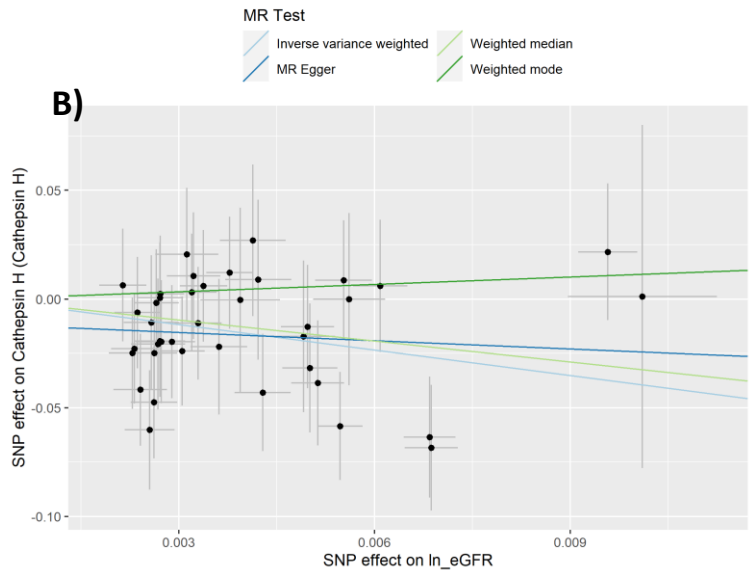
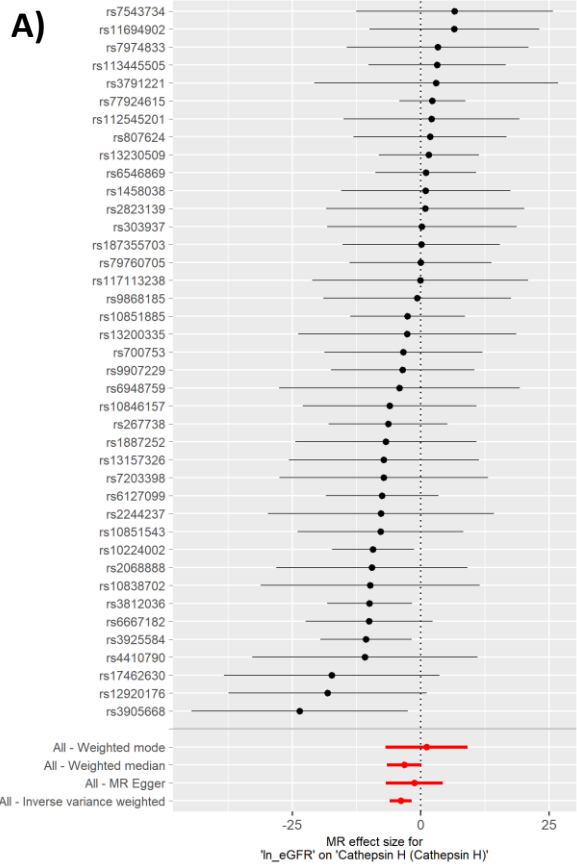


Suppl. Fig. 13. Forward MR analysis for eGFR-CST6

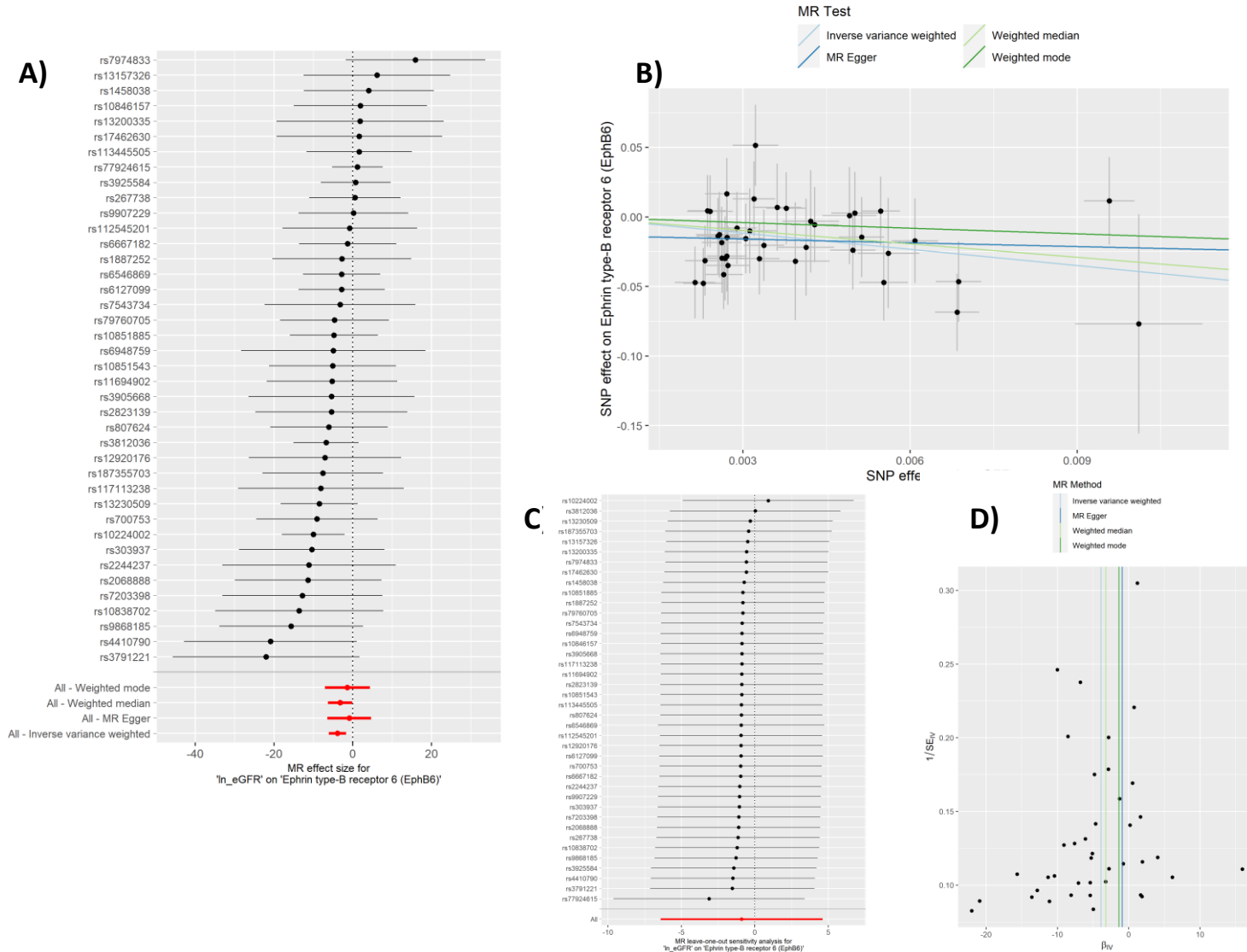


Suppl. Fig. 14. Forward MR analysis for eGFR-Cathepsin H

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41

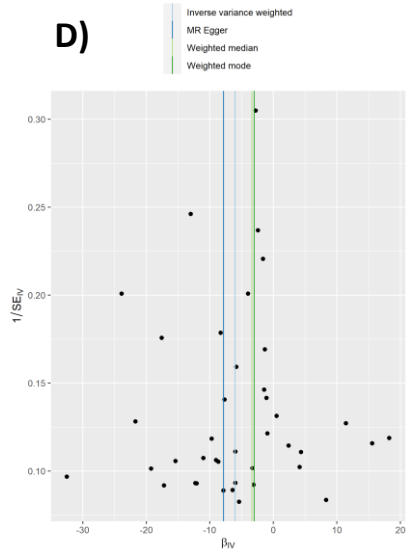
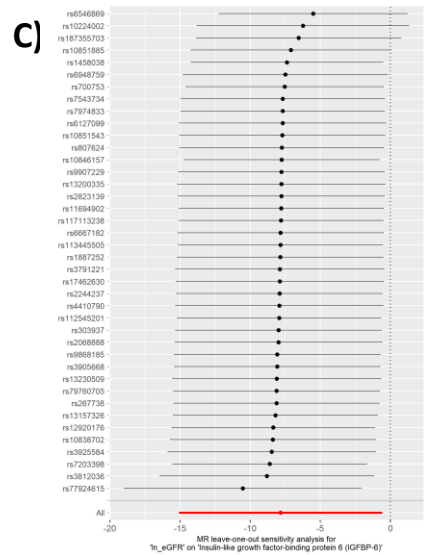
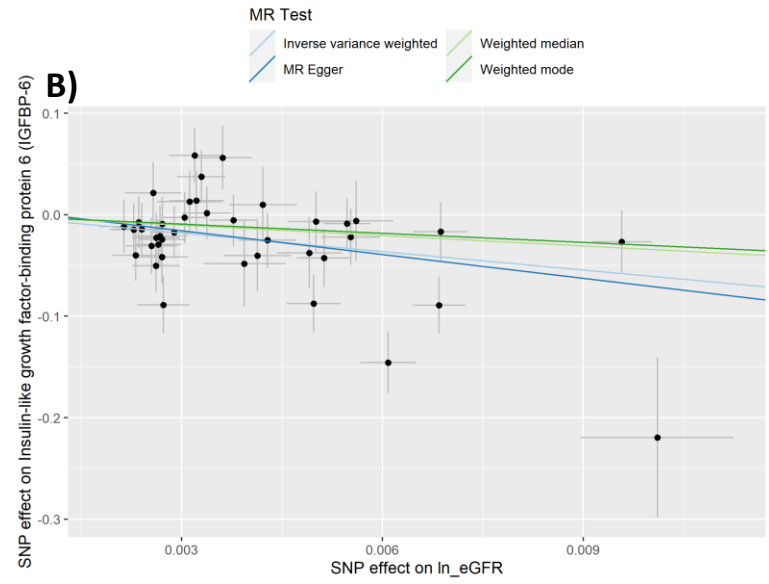
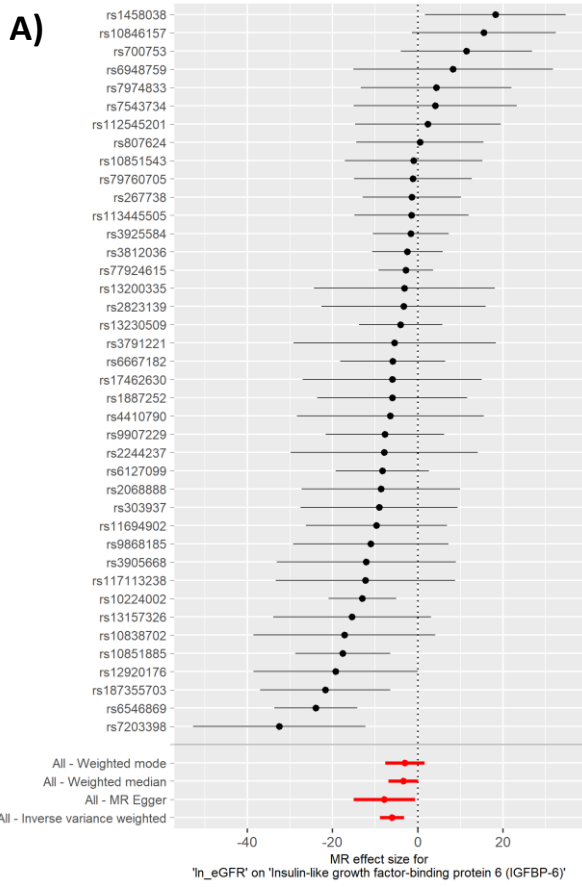


Suppl. Fig. 15. Forward MR analysis for eGFR-EphB6

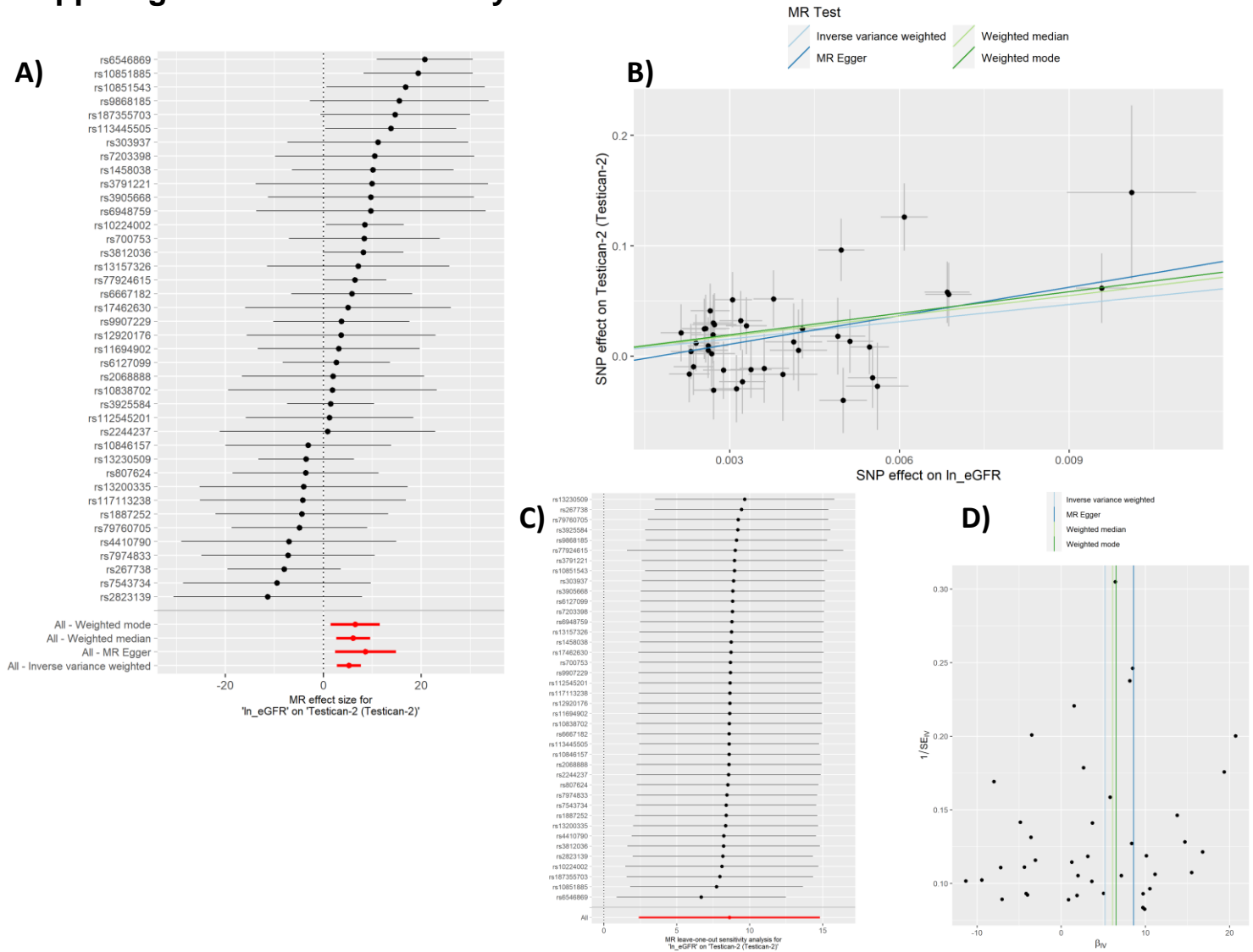


Suppl. Fig. 16. Forward MR analysis for eGFR-IGFBP-6

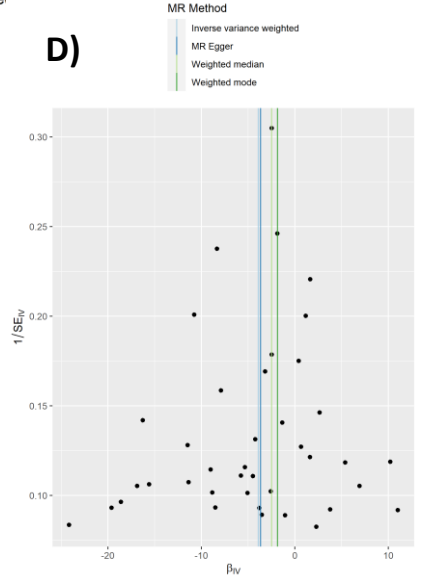
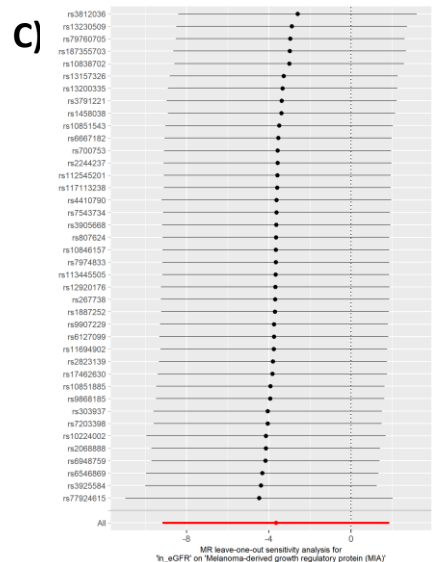
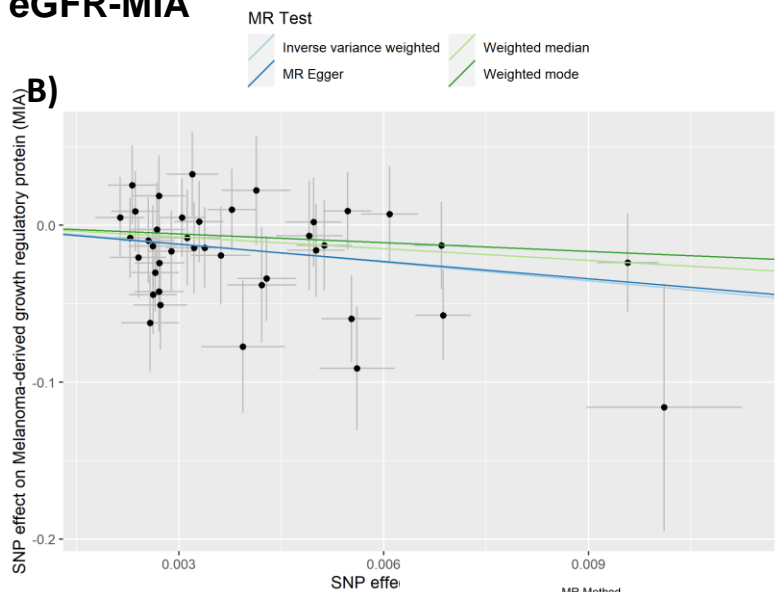
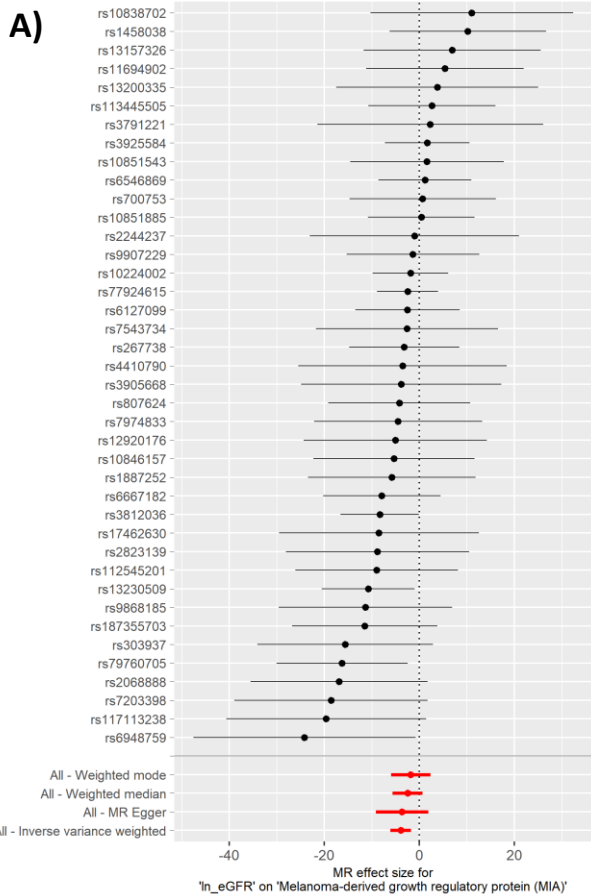
1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41



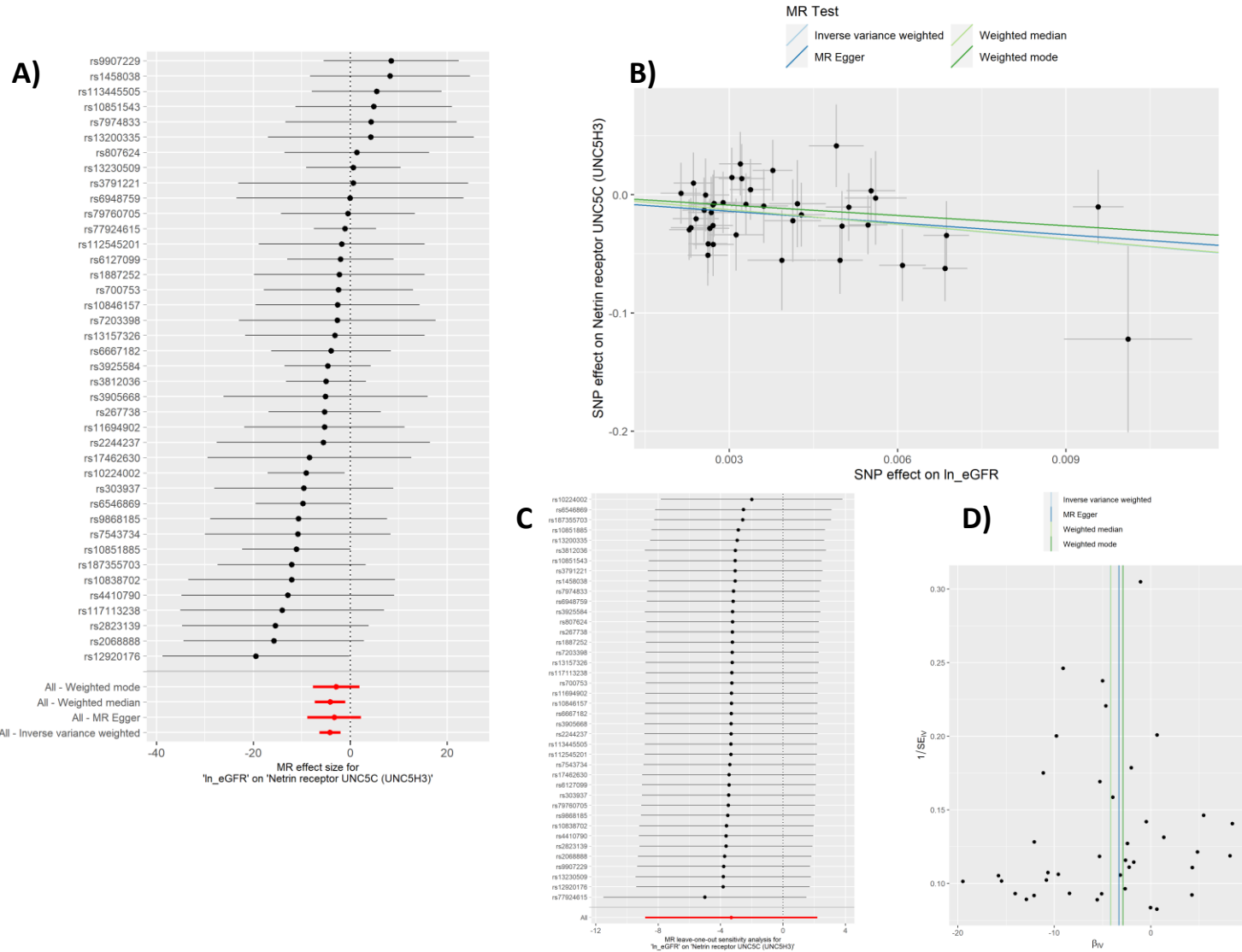
Suppl. Fig. 17. Forward MR analysis for eGFR-SPOCK2



Suppl. Fig. 18. Forward MR analysis for eGFR-MIA

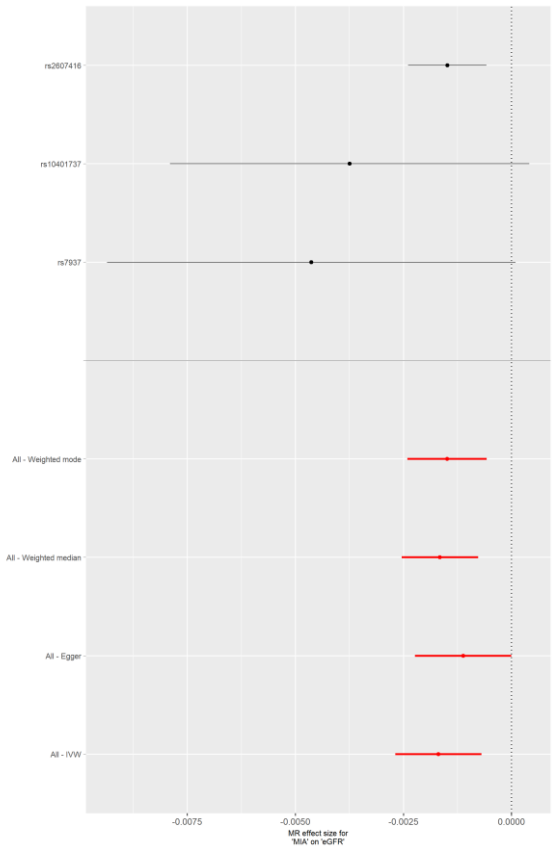


Suppl. Fig. 19. Forward MR analysis for eGFR-UNC5H3



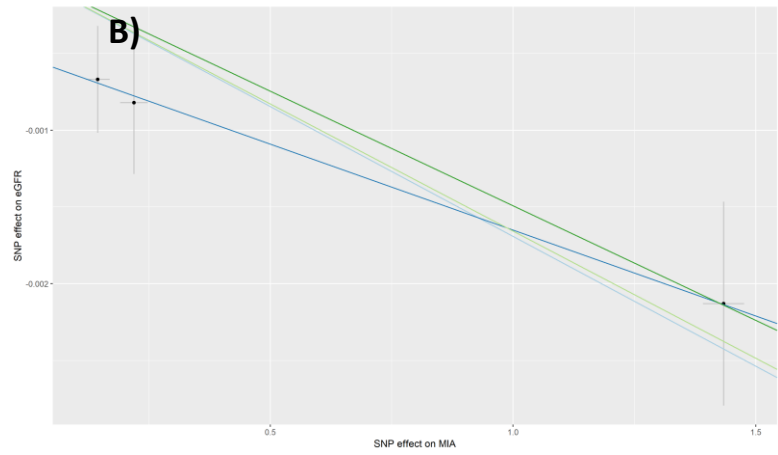
Suppl. Fig. 20. Reverse MR analysis for MIA-eGFR

A)

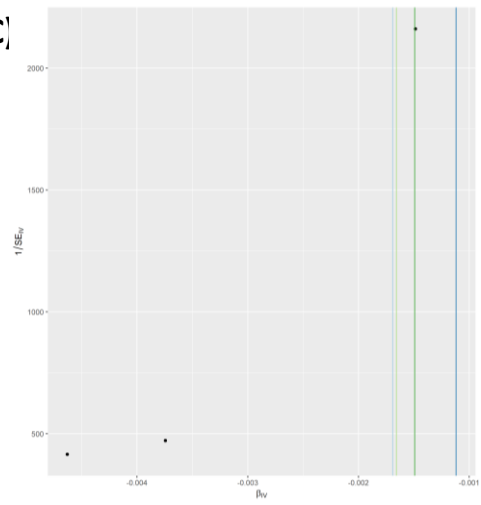


MR Test
 Inverse variance weighted
 MR Egger
 Weighted median
 Weighted mode

B)

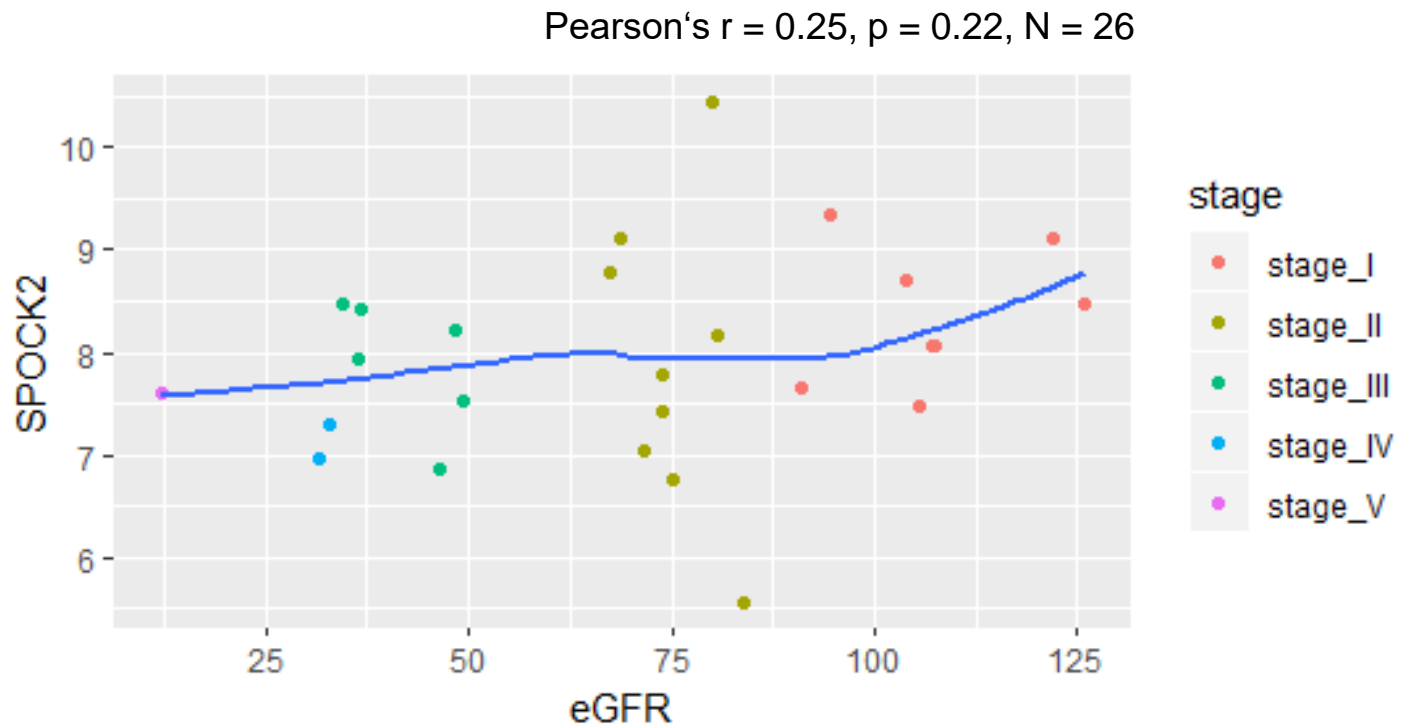


C)



1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41

Suppl. Fig. 21. *SPOCK2* gene expression in renal tissue from 26 CKD patients



Manuscript ID:

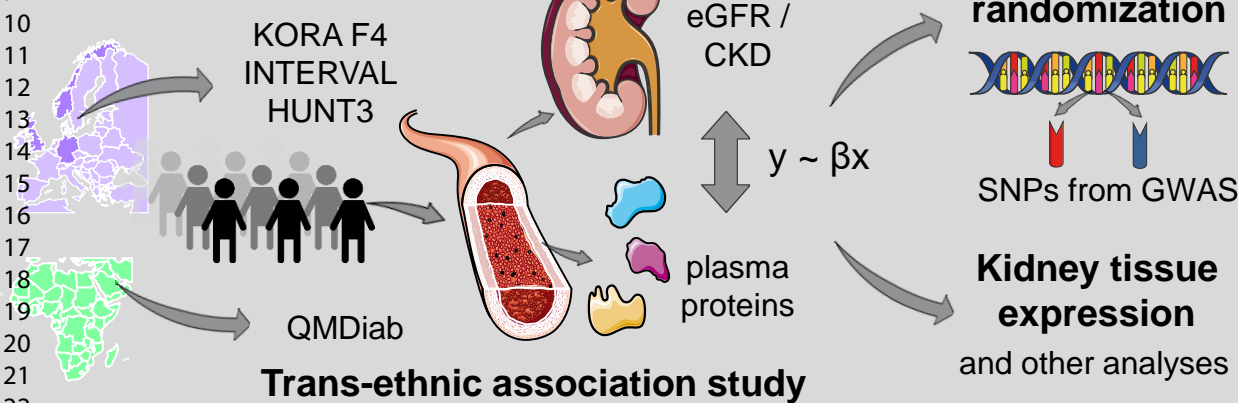
Adapted* STROBE Checklist for Observational Studies

Section/Topic	Checklist Item	Response
Background	Specific objectives clearly stated, including any pre-specified hypotheses.	
METHODS		
Study Design	Key elements of study design presented early in the paper.	
Setting	Setting and location from where study participants drawn clearly described.	
	Relevant dates for periods of recruitment, exposure, follow-up, and data collection.	
Participants	Eligibility criteria or matching criteria, as appropriate, clearly described.	
	Sources and methods of selection of participant (or cases and controls) selection.	
	Number of exposed/unexposed OR number of controls per case, as appropriate.	
	Methods of follow-up clearly described.	
Variables	Exposures, outcomes, predictors, potential confounders, effect modifiers clearly defined.	
Data Sources/ Measurement	Sources of data and details of methods of assessment.	
	Comparability of assessment methods, if more than one group.	
Bias	Clear description of efforts to address potential sources of bias.	
Statistical Methods	All statistical methods clearly described, including those used to control for confounding.	
	Methods used to examine subgroups and interactions clearly described.	
	Extent of missing data, and how it was handled clearly described.	
	Extent of loss to follow-up and how it was addressed clearly described.	
	Sensitivity Analyses appropriate and clearly described.	
RESULTS		
Participants	Number of individuals at each stage – e.g., number eligible, examined for eligibility, confirmed eligible, included in study, completing follow-up and analyzed.	
	Reasons for non-participation at each stage clearly described.	
Descriptive Data	Characteristics of study participants presented .	
	Information on exposures and potential confounders available.	
	Number of participants with missing data for each variable of interest provided.	
	Follow-up time summarized (average and total, as appropriate).	
Outcomes	Number of outcome events or summary measures over time available.	
	For case-control study, numbers in each exposure category, or summary measures of exposure provided.	
	For case-control study, numbers of outcomes events or summary measures.	
Main Results	Unadjusted estimates and, if applicable, confounder-adjusted estimates and their precision (95% confidence interval) provided.	
	Category boundaries, when continuous variables were categorized, available.	
	Estimates of relative risk translated into absolute risk for a meaningful period.	
Other Analyses	Sub-group analyses, interactions, sensitivity analyses adequately presented.	

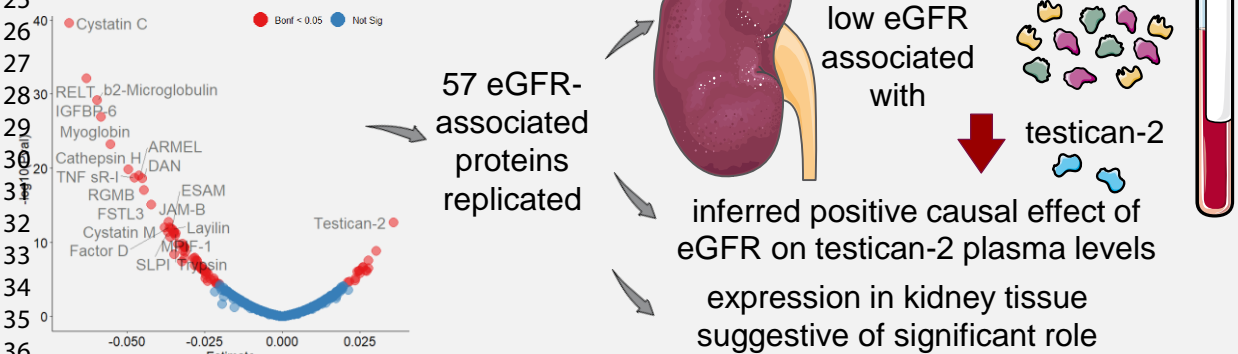
Additional details:

Plasma proteomics of renal function: a trans-ethnic meta-analysis and Mendelian randomization study

Methods



Outcome



Conclusion

In the largest multiplex plasma proteomics study to date, we identify 57 proteins as trans-ethnically associated with eGFR and/or CKD. Mendelian randomization analysis suggests that several proteins have a causal relationship with kidney function, results which highlight in particular testican-2. This work represents an early milestone in the identification and establishment of sets of proteins that could act as physiological markers of kidney disease progression, biomarkers of potential clinical relevance.