

# Global dataset of soil eukaryotic communities created with a uniform protocol and long read sequencing

Received: 8 December 2025

Accepted: 22 April 2026

Cite this article as: Mikryukov, V., Dulya, O., Abarenkov, K. *et al.* Global dataset of soil eukaryotic communities created with a uniform protocol and long read sequencing. *Sci Data* (2026). <https://doi.org/10.1038/s41597-026-07315-y>

Vladimir Mikryukov, Olesya Dulya, Kessy Abarenkov, Sten Anslan, Niloufar Hagh-Doust, Victoria Prins, Kristel Panksep, Sergei Pölme, Khalid S. Ibrahim, Mo Bahram, Kalev Adamson, Ahto Agan, Talaat Ahmed, Juha M. Alatalo, Felipe E. Albornoz, Abdullah MS Al-Hatmi, Saad Alkahtani, Julieta Alvarez-Manjarrez, Jelena Ankuda, Alexandre Antonelli, Manikandan Ariyan, Kęstutis Armolaitis, Farzad Aslani, Isabel C. Barrio, Marijn Bauters, Elisabeth Machteld Biersma, Krišs Bitenieks, Gregory Bonito, Francis Q. Brearley, Kari Anne Bråthen, Franz Buegger, Klaus Butterbach-Bahl, Miklós Bálint, Erin K. Cameron, Fabiana Canini, Rebeca Casique-Valdés, Adriana Corrales, Evgeny A. Davydov, Eske De Crop, André De Kesel, Joseph Fovo Djeugap, Rein Drenkhan, Camila Duarte Ritter, Sergey V. Dudov, Mikk Espenberg, Ongua Fanuel, Vladimir E. Fedosov, Luke Florence, Brendan R. Furneaux, Ariadne N. M. Furtado, Sanni Färkkilä, Natalia S. Gamova, Roberto Garibay-Orijel, József Geml, Soumya Ghosh, Roberto Godoy, Daniyal Gohar, Marieka Gryzenhout, Ayad H. Hasan, Amr H. Hashem, Jacob Heilmann-Clausen, Terry W. Henkel, Indrek Hiiesalu, Inga Hiiesalu, Mahdieh S. Hosseyni Moghaddam, Kevin D. Hyde, Karina Inostroza, Khalil Kariman, Elina Karimullina, Sebastian Kepfer-Rojas, Abdul Nasir Khalid, Darta Klavina, Petr Kohout, Yuri N. Korotkov, John Y. Kupagme, Olavi Kurina, Louis James Lamit, Adebola Azeez Lateef, Njouonkou André Ledoux, Young Woon Lim, Jose G. Maciá-Vicente, Kristaps Makovskis, Sebastián Martínez, César Marín, Peter Meidl, Peter E. Mortimer, Sunil Mundra, Victoria Naluyange, Tarquin Netherway, Kevin K. Newsham, Eduardo Nouhra, Casper Nyamukondiwa, Vincent Nteziryayo, Dennis M. W. Ochieno, Jane Oja, Vladimir G. Onipchenko, Eveli Otsing, Mustafa Nadhim Owaid, Meike Piepenbring, Polina Pochekutova, Maihyra Marina Pombo, Karin Pritsch, Rasmus Puusepp, Jaan Pärn, Kadri Pöldmaa, Saleh Rahimlou, Andrea C. Rinaldi, Oscar Rojas, Tomas Roslin, Kadri Runnel, Elisabeth Rähn, Malka Saba, Alessandro Saitta, Talal Sabhan Salih, Joosep Sarapuu, Eduard Serrano, Oscar Serrano, Dipon Sharmah, Cathy Sharp, Maria W. Skalska-Tuomi, Kassim Issifou Tchan, Camille Truong, Helga van der Merwe, Linda L. P. Vanié-Léabo, Aida M. Vasco-Palacios, Annemieke Verbeken, Lukáš Vlk, Nalin N. Wijayawardene, Jennifer L. Wood, W. A. Erandi Yasanthika, Nourou S. Yorou, Geoffrey Zahn, Irma Zettur, Laura Zucconi, Urmas Kõljalg & Leho Tedersoo

We are providing an unedited version of this manuscript to give early access to its findings. Before final publication, the manuscript will undergo further editing. Please note there may be errors present which affect the content, and all legal disclaimers apply.

If this paper is publishing under a Transparent Peer Review model then Peer Review reports will publish with the final article.

ARTICLE IN PRESS

## Global dataset of soil eukaryotic communities created with a uniform protocol and long read sequencing

Vladimir Mikryukov 0000-0003-2786-2690<sup>1,2</sup>✉, Olesya Dulya 0000-0002-7185-8659<sup>1,2</sup>, Kessy Abarenkov 0000-0001-5526-4845<sup>3</sup>, Sten Anslan 0000-0002-2299-454X<sup>4,2</sup>, Niloufar Hagh-Doust 0000-0003-0616-5829<sup>1,2</sup>, Victoria Prins 0009-0003-8968-6773<sup>5,2</sup>, Kristel Panksep 0000-0003-4743-6111<sup>6,7,5</sup>, Sergei Pölme 0000-0002-9658-1166<sup>8,2</sup>, Khalid S. Ibrahim 0000-0002-8585-3429<sup>9</sup>, Mo Bahram 0000-0002-9539-3307<sup>10,11</sup>, Kalev Adamson 0000-0002-8810-8838<sup>12</sup>, Ahto Agan<sup>12</sup>, Talaat Ahmed 0000-0001-8022-1855<sup>13</sup>, Juha M. Alatalo 0000-0001-5084-850X<sup>13</sup>, Felipe E. Alborno 0000-0001-9526-0945<sup>14</sup>, Abdullah MS Al-Hatmi 0000-0002-5206-2647<sup>15</sup>, Saad Alkahtani 0000-0001-7381-5110<sup>16</sup>, Julieta Alvarez-Manjarrez 0000-0002-5581-7443<sup>17</sup>, Jelena Ankuda 0000-0002-9446-7802<sup>18</sup>, Alexandre Antonelli 0000-0003-1842-9297<sup>19,20,21</sup>, Manikandan Ariyan 0000-0002-7277-0544<sup>1</sup>, Kęstutis Armolaitis 0000 0001 8295 2440<sup>22</sup>, Farzad Aslani 0000-0003-0485-9800<sup>23</sup>, Isabel C. Barrio 0000-0002-8120-5248<sup>24</sup>, Marijn Bauters 0000-0003-0978-6639<sup>25</sup>, Elisabeth Machteld Biersma 0000-0002-9877-2177<sup>26</sup>, Krišs Biteniekš 0000-0001-8829-3263<sup>27</sup>, Gregory Bonito<sup>28</sup>, Francis Q. Brearley 0000-0001-5053-5693<sup>29</sup>, Kari Anne Bråthen 0000-0003-0942-1074<sup>30</sup>, Franz Buegger 0000-0003-3526-4711<sup>31</sup>, Klaus Butterbach-Bahl 0000-0001-9499-6598<sup>32</sup>, Miklós Bálint 0000-0003-0499-8536<sup>33,34,35</sup>, Erin K. Cameron 0000-0002-3374-9830<sup>36</sup>, Fabiana Canini 0000-0002-9626-6318<sup>37</sup>, Rebeca Casique-Valdés 0000-0002-0497-8255<sup>38</sup>, Adriana Corrales 0000-0001-9885-4634<sup>39</sup>, Evgeny A Davydov 0000-0002-2316-8506<sup>40</sup>, Eske De Crop 0000-0002-9067-6981<sup>41</sup>, André De Kesel 0000-0002-5287-4637<sup>42</sup>, Joseph Fovo Djeugap 0000-0003-4786-7026<sup>43</sup>, Rein Drenkhan<sup>12</sup>, Camila Duarte Ritter 0000-0002-3371-7425<sup>44</sup>, Sergey V. Dudov 0000-0003-1512-0956<sup>40</sup>, Mikk Espenberg 0000-0003-0469-6394<sup>45</sup>, Ongua Fanuel 0009-0002-4902-4858<sup>46</sup>, Vladimir E. Fedosov 0000-0002-5331-6346<sup>40</sup>, Luke Florence 0000-0002-1901-7772<sup>47</sup>, Brendan R. Furneaux 0000-0003-3522-7363<sup>4</sup>, Ariadne N.M. Furtado 0000-0001-9596-7553<sup>48</sup>, Sanni Färkkilä 0000-0002-4485-5567<sup>49</sup>, Natalia S. Gamova 0000-0002-4141-757X<sup>40</sup>, Roberto Garibay-Orijel 0000-0002-6977-7550<sup>50</sup>, József Geml 0000-0001-8745-0423<sup>51</sup>, Soumya Ghosh 0000-0002-4945-3516<sup>15,52</sup>, Roberto Godoy 0000-0002-3719-3091<sup>53</sup>, Daniyal Gohar 0000-0003-0312-1142<sup>54</sup>, Marieka Gryzenhout<sup>55</sup>, Ayad H. Hasan 0000-0001-7280-2254<sup>56</sup>, Amr H. Hashem<sup>57</sup>, Jacob Heilmann-Clausen 0000-0003-4713-6004<sup>58</sup>, Terry W. Henkel 0000-0001-9760-8837<sup>59</sup>, Indrek Hiiesalu<sup>60</sup>, Inga Hiiesalu 0000-0002-5457-2376<sup>1</sup>, Mahdieh Hosseyni Moghadam<sup>1</sup>, Kevin D. Hyde<sup>61,62</sup>, Karina Inostroza 0000-0003-3026-6291<sup>63</sup>, Khalil Kariman 0000-0002-2070-4713<sup>64</sup>, Elina Karimullina 0000-0001-6459-2293<sup>65</sup>, Sebastian Kepfer-Rojas 0000-0002-1681-2877<sup>66</sup>, Abdul Nasir Khalid<sup>67</sup>, Darta Klavina 0000-0002-1455-9062<sup>27</sup>, Petr Kohout 0000-0002-3985-2310<sup>68</sup>, Yuri N. Korotkov 0000-0001-8436-0887<sup>40</sup>, John Y. Kupagme 0000-0002-9981-050X<sup>2</sup>, Olavi Kurina 0000-0002-4858-4629<sup>69</sup>, Louis James Lamit 0000-0002-0385-6010<sup>70</sup>, Adebola Azeez Lateef<sup>71,72</sup>, Njouonkou André Ledoux 0000-0001-9733-7248<sup>73</sup>, Young Woon Lim 0000-0003-2864-3449<sup>74</sup>, Jose G. Maciá-Vicente<sup>75</sup>, Kristaps Makovskis 0000-0003-4943-1912<sup>27</sup>, Sebastián Martínez 0000-0003-0455-5823<sup>76</sup>, César Marín 0000-0002-2529-8929<sup>77,78</sup>, Peter Meidl<sup>79</sup>, Peter E. Mortimer<sup>80,81</sup>, Sunil Mundra 0000-0002-0535-118X<sup>82,83</sup>, Victoria Naluyange 0000-0002-3990-9751<sup>84</sup>, Tarquin Netherway 0000-0002-9049-9225<sup>10</sup>, Kevin K. Newsham 0000-0002-9108-0936<sup>85</sup>, Eduardo Nouhra 0000-0002-7080-8211<sup>86</sup>, Casper Nyamukondiwa<sup>87,88</sup>, Vincent Ntezirayayo 0000-0002-0176-5780<sup>89</sup>, Dennis M.W. Ochieno 0000-

0003-3985-9421<sup>90</sup>, Jane Oja 0000-0003-1446-2284<sup>1</sup>, Vladimir G. Onipchenko 0000-0002-1626-1171<sup>91</sup>, Eveli Otsing 0000-0001-7416-257X<sup>2</sup>, Mustafa Nadhim Owaid 0000-0001-9005-4368<sup>92</sup>, Meike Piepenbring 0000-0002-7043-5769<sup>93</sup>, Polina Pochekutova 0009-0001-0453-9869<sup>2</sup>, Maihyra Marina Pombo 0000-0002-0329-9736<sup>94</sup>, Karin Pritsch 0000-0001-6384-2473<sup>31</sup>, Rasmus Puusepp 0000-0002-0617-3776<sup>2</sup>, Jaan Pärn 0000-0001-6507-8894<sup>45</sup>, Kadri Põldmaa 0000-0002-7936-2455<sup>1</sup>, Saleh Rahimlou 0000-0003-0427-1329<sup>95</sup>, Andrea C. Rinaldi 0000-0002-9352-1037<sup>48</sup>, Oscar Rojas<sup>96</sup>, Tomas Roslin 0000-0002-2957-4791<sup>11,97</sup>, Kadri Runnel 0000-0002-7308-3623<sup>98</sup>, Elisabeth Rähn<sup>12</sup>, Malka Saba 0000-0001-7673-2345<sup>99</sup>, Alessandro Saitta<sup>100</sup>, Talal Sabhan Salih 0000-0001-8345-6131<sup>101</sup>, Joosep Sarapuu 0000-0002-4805-5612<sup>102</sup>, Eduard Serrano 0000-0001-9082-473X<sup>103,63</sup>, Oscar Serrano 0000-0002-5973-0046<sup>103</sup>, Dipon Sharmah 0000-0002-5579-3103<sup>104</sup>, Cathy Sharp 0009-0003-4985-1543<sup>105</sup>, Maria W. Skalska-Tuomi 0000-0002-7154-5177<sup>106,107</sup>, Kassim Tchan Issifou 0000-0002-8572-3405<sup>108</sup>, Camille Truong 0000-0002-8510-1761<sup>109,110</sup>, Helga van der Merwe 0000-0002-7677-5123<sup>111,112</sup>, Linda L.P. Vanié-Léabo 0000-0001-6504-3456<sup>113</sup>, Aida M. Vasco-Palacios 0000-0003-0539-9711<sup>114</sup>, Annemieke Verbeke<sup>41</sup>, Lukáš Vlk 0000-0001-9201-8715<sup>115</sup>, Nalin N. Wijayawardene 0000-0003-0522-5498<sup>116</sup>, Jennifer L. Wood 0000-0002-7313-5681<sup>117</sup>, W.A. Erandi Yasanthika 0000-0002-3757-3801<sup>62</sup>, Nourou S. Yorou 0000-0001-6997-811X<sup>118</sup>, Geoffrey Zahn 0000-0002-8691-692X<sup>119</sup>, Irma Zettur<sup>3</sup>, Laura Zucconi 0000-0001-9793-2303<sup>37</sup>, Urmas Kõljalg 0000-0002-5171-1668<sup>1</sup> & Leho Tedersoo 0000-0002-1635-1249<sup>120,16</sup>✉

<sup>1</sup>Department of Botany, Institute of Ecology and Earth Sciences, University of Tartu, Juhan Liivi 2, Tartu, 50409, Estonia, <sup>2</sup>Mycology and Microbiology Center, University of Tartu, Juhan Liivi 2, Tartu, 50409, Estonia, <sup>3</sup>Natural History Museum, University of Tartu, Vanemuise 46, Tartu, 51003, Estonia, <sup>4</sup>Department of Biological and Environmental Science, Faculty of Mathematics and Science, University of Jyväskylä, Survantie 9C, Jyväskylä, Central Finland, 40500, Finland, <sup>5</sup>Institute of Technology, Faculty of Science and Technology, University of Tartu, Nooruse 1, Tartu, 50411, Estonia, <sup>6</sup>Department of hydrobiology and fisheries, Institute of Agricultural and Environmental Sciences, Estonian University of Life Sciences, F. R. Kreutzwaldi 1, Tartu, 51006, Estonia, <sup>7</sup>Institute of Freshwater Research, Department of Aquatic Resources, Swedish University of Agricultural Sciences, Stångholmsvägen 2, Drottningholm, 178 93, Sweden, <sup>8</sup>Natural History Museum and Botanical Garden, University of Tartu, Vanemuise 46, Tartu, 50410, Estonia, <sup>9</sup>Department of Biology, College of Science, University of Zakho, Zakho, Kurdistan Region, 42002, Iraq, <sup>10</sup>Department of Agroecology, Aarhus University, Forsøgsvej 1, Slagelse, 4200, Denmark, <sup>11</sup>Department of Ecology, Swedish University of Agricultural Sciences (SLU), Ulls väg 18B, Uppsala, 75651, Sweden, <sup>12</sup>Chair of Silviculture and Forest Ecology, Institute of Forestry and Engineering, Estonian University of Life Sciences, F. R. Kreutzwaldi 5, Tartu, 51006, Estonia, <sup>13</sup>Environmental Science Center, Qatar University, Doha, Qatar, <sup>14</sup>School of Biological Sciences, The University of Western Australia, 35 Stirling Hwy, Crawley, WA, 6009, Australia, <sup>15</sup>Natural and Medical Sciences Research Center, University of Nizwa, Nizwa, 616, Oman, <sup>16</sup>Department of Zoology, College of Science, King Saud University, P.O. Box 2455, Riyadh, 11451, Saudi Arabia, <sup>17</sup>Micología Integral, Instituto de Biología, Universidad Nacional Autónoma de México, Tercer circuito interior s/n, Ciudad Universitaria, Coyoacán, Mexico City, 04510, México, <sup>18</sup>Vokė Branch, Institute of Agriculture, Lithuanian Research Centre for Agriculture and Forestry, Žalioji Sq. 2, Vilnius, LT-02232, Lithuania, <sup>19</sup>Royal Botanic Gardens, Kew GB, Richmond, Surrey, TW9 3AE, United Kingdom, <sup>20</sup>Gothenburg Global Biodiversity Centre, Department of Biological and Environmental Sciences, University of Gothenburg, Göteborg, 405 30, Sweden, <sup>21</sup>Department of Biology, University of Oxford, South Parks Road, Oxford, OX1 3RB, United Kingdom, <sup>22</sup>Department of Silviculture and Ecology, Institute of Forestry, Lithuanian Research Centre for Agriculture and Forestry, Liepų str. 1, Girionys, Kaunas distr., LT-53101, Lithuania, <sup>23</sup>Above-belowground interactions group, Institute of Biology, Leiden University, Rapenburg 70, Leiden, Netherlands, <sup>24</sup>Faculty of Life and Land, Agricultural University of Iceland, Árleyni 22, Reykjavík, Iceland, 112, Iceland, <sup>25</sup>Q-ForestLab, Department of Environment, Ghent University, Coupure Links 653,

Gent, 9000, Belgium, <sup>26</sup>Natural History Museum of Denmark, University of Copenhagen, Gothersgade 130, Copenhagen, 1123, Denmark, <sup>27</sup>Latvian State Forest Research Institute "Silava", Rīga str. 111, Salaspils, LV-2169, Latvia, <sup>28</sup>Department of Plant Soil and Microbial Sciences, Michigan State University, Michigan, MI, 48824, USA, <sup>29</sup>Department of Natural Sciences, Manchester Metropolitan University, Chester Street, Manchester, M1 5GD, UK, <sup>30</sup>Department of Arctic and Marine Ecology, UiT The Arctic University of Norway, Breivika, Tromsø, 9037, Norway, <sup>31</sup>Research Unit Environmental Simulation, Helmholtz Zentrum München, Ingolstädter Landstraße 1, Neuherberg, 85764, Germany, <sup>32</sup>Agroecology, Department, Aarhus University, Ole Worms Alle 3, Aarhus, 8000, Denmark, <sup>33</sup>Functional Environmental Genomics, Senckenberg Biodiversity and Climate Research Centre, Georg-Voigt-Str. 14-16, Frankfurt am Main, Hesse, 60325, Germany, <sup>34</sup>Institute of Insect Biotechnology, Justus Liebig University, Heinrich-Buff-Ring 26, Gießen, Hesse, 35392, Germany, <sup>35</sup>Functional Environmental Genomics, LOEWE Centre for Translational Biodiversity Genomics, Senckenberganlage 25, Frankfurt am Main, Hesse, 60325, Germany, <sup>36</sup>Department of Biology, Saint Mary's University, 923 Robie St., Halifax, Nova Scotia, B3H 2M4, Canada, <sup>37</sup>Department of Ecological and Biological Sciences, University of Tuscia, Largo dell'Università, Viterbo, 01100, Italy, <sup>38</sup>Horticulture department, Universidad Autónoma Agraria Antonio Narro, Saltillo, Coahuila, 25315, México, <sup>39</sup>Society for the Protection of Underground Networks (SPUN), 500 South DuPont Highway, Dover, DE, 19901, USA, <sup>40</sup>no affiliation, Russia, <sup>41</sup>Department of Biology, Ghent University, K.L. Ledeganckstraat 35, Gent, 9000, Belgium, <sup>42</sup>Research Department, Biodiversity and Evolution, Meise Botanic Garden, Nieuwelaan 38, Meise, B-1860, Belgium, <sup>43</sup>Department of Crop Sciences, Plant Protection, University of Dschang, Dschang, West Region, PO. Box 222, Cameroon, <sup>44</sup>Intituto Juruá, Manaus, Amazonas, 69083-300, Brazil, <sup>45</sup>Department of Geography, Institute of Ecology and Earth Sciences, University of Tartu, Vanemuise 46, Tartu, 51003, Estonia, <sup>46</sup>Soils, Environment and Agrometeorology, National Agricultural Research Laboratories-Kawanda, National Agricultural Research Organization NARO, Kampala, 7065, Uganda, <sup>47</sup>Department of Ecological, Plant and Animal Sciences, La Trobe University, Bundoora, Victoria, 3086, Australia, <sup>48</sup>Department of Biomedical Sciences, University of Cagliari, Cittadella Universitaria, Monserrato (CA), I09042, Italy, <sup>49</sup>University of Tartu, Ülikooli 18a, Tartu, 51005, Estonia, <sup>50</sup>Universidad Nacional Autónoma de México, Instituto de Biología, Circuito exterior s/n, Ciudad de México, Ciudad de México, 4510, México, <sup>51</sup>Environmental Microbiome Research Group, Research and Development Centre, Eszterházy Károly Catholic University, Leányka u. 8., Eger, 3300, Hungary, <sup>52</sup>University of the Free State, Bloemfontein, 9301, South Africa, <sup>53</sup>Instituto de Ciencias Ambientales y Evolutivas, Universidad Austral de Chile, Isla Teja sn., Valdivia, 5090000, Chile, <sup>54</sup>Department of Botany and Plant Pathology, Oregon State University, 2701 SW Campus Way, Corvallis, OR, 97331, USA, <sup>55</sup>Department of Genetics, Natural and Agricultural Sciences, University of the Free State, Bloemfontein, 9300, South Africa, <sup>56</sup>Department of Medical Microbiology, Faculty of Science and Health, Koya University, Koya 44023, Kurdistan Region – F.R., Iraq, <sup>57</sup>Botany and Microbiology Department, Faculty of Science, Al-Azhar University, Nasr City, Cairo, 11884, Egypt, <sup>58</sup>Center for Macroecology, Evolution and Climate, Globe institute, University of Copenhagen, Universitetsparken 15, Copenhagen, 2100, Denmark, <sup>59</sup>Department of Biological Sciences, California State Polytechnic University, Humboldt, Arcata, CA, 95521, USA, <sup>60</sup>Institute of Ecology and Earth Sciences, University of Tartu, Juhan Liivi 2, Tartu, 50409, Estonia, <sup>61</sup>Department of Botany and Microbiology, College of Science, King Saud University, P.O. Box 2455, Riyadh, 11495, Saudi Arabia, <sup>62</sup>Center of Excellence in Fungal Research, Mae Fah Luang University, Chiang Rai, 57100, Thailand, <sup>63</sup>BIOSFERA Research & Conservation, DS Can Mutxo, Sils, Girona, 17410, Spain, <sup>64</sup>UWA School of Agriculture and Environment, The University of Western Australia, Perth, WA, 6009, Australia, <sup>65</sup>Microbiology, Immunology and Infectious Diseases, University of Calgary, 3330 Hospital Drive NW, Calgary, Alberta, T2N 4N1, Canada, <sup>66</sup>Geosciences and Natural Resource Management, University of Copenhagen, Rolighedsvej 23, Frederiksberg, 1958, Denmark, <sup>67</sup>Institute of Botany, University of the Punjab, Lahore, Punjab, 54320, Pakistan, <sup>68</sup>Laboratory of Microbial Ecology and Biogeography, Institute of Microbiology, Czech Academy of Science, Videnska 1083, Prague, 14220, Czechia, <sup>69</sup>Chair of Biological Diversity and Nature Tourism, Institute of Agricultural and Environmental Sciences, Estonian University of Life Sciences, F. R. Kreutzwaldi 5, Tartu, 51006, Estonia, <sup>70</sup>Department of Biology, Syracuse University, 107 College Place, Syracuse, NY, 13244, USA, <sup>71</sup>Department of Plant Biology, Faculty of Life Sciences, University of Ilorin, Ilorin, Kwara, 1515, Nigeria, <sup>72</sup>Department of Forest Sciences, Faculty of Forestry and Agriculture, University of Helsinki, Agnes Sjöbergin katu 2, P.O. Box 66, Helsinki, 14, Finland, <sup>73</sup>Department of Plant Sciences, Faculty of Science, The University of Bamenda, Bambili, North-West Region, PO Box 39, Cameroon, <sup>74</sup>School of Biological Sciences, College of Natural Science, Institute of Biodiversity, Gwanak-ro 1, Seoul, 8826, South Korea, <sup>75</sup>Department of Marine Sciences and Applied Biology, University of Alicante, Alicante, 3080, Spain, <sup>76</sup>Laboratorio de Patología Vegetal, Instituto Nacional de Investigación Agropecuaria, Ruta 8, Km 281, Treinta y Tres,

33000, Uruguay, <sup>77</sup>Centro de Investigación e Innovación para el Cambio Climático (CiiCC), Universidad Santo Tomás, Av. Ramón Picarte 1130, Valdivia, 5090000, Chile, <sup>78</sup>Amsterdam Institute for Life and Environment, Section Ecology & Evolution, Vrije Universiteit Amsterdam, de Boelelaan 1085, Amsterdam, 1081 HV, Netherlands, <sup>79</sup>Soil Ecology, Free University Berlin, Berlin, 13357, Germany, <sup>80</sup>Soil Science, Stellenbosch University, PO BOX X1, Matieland, Stellenbosch, Western Cape, 7600, South Africa, <sup>81</sup>Applied Symbiotics, 111 9th Road, Hyde Park, Gauteng, 2196, South Africa, <sup>82</sup>Department of Biology, College of Science, United Arab Emirates University, Al Ain, Abu Dhabi, UAE, <sup>83</sup>Khalifa Center for Genetic Engineering and Biotechnology, United Arab Emirates University, Al Ain, Abu Dhabi, UAE, <sup>84</sup>Department of Land Use Management, School of Agriculture, Veterinary Sciences and Technology, Masinde Muliro University of Science and Technology, Kakamega, 190-50100, Kenya, <sup>85</sup>British Antarctic Survey, Natural Environment Research Council, Madingley Road, Cambridge, CB3 0ET, United Kingdom, <sup>86</sup>CONICET-FCEfyN, National University of Córdoba, Av. Vélez Sarsfield 1666, Córdoba, Córdoba, X5016GCN, Argentina, <sup>87</sup>Department of Biological Sciences and Biotechnology, Botswana International University of Science and Technology, P. Bag 16, Palapye, Botswana, <sup>88</sup>Centre for Environmental Policy, Imperial College London, Kennedy Building, Silwood Park, Ascot, Berkshire, SL5 7PY, United Kingdom, <sup>89</sup>Food Science and Technology, Faculty of Agriculture and Bio-engineering, University of Burundi, Avenue l'Unesco no 2, Bujumbura, Burundi, 2940, Burundi, <sup>90</sup>Department of Biological Sciences, School of Natural Sciences, Masinde Muliro University of Science and Technology, Kakamega-Webuye Road, Kakamega, 254, Kenya, <sup>91</sup>Biological Faculty, Shenzhen MSU-BIT University, Shenzhen, 518115, China, <sup>92</sup>Department of Environment, College of Applied Sciences-Hit, University Of Anbar, Hit, Anbar, Iraq, <sup>93</sup>Mycology, Goethe University, Max-von-Laue-Str. 13, Frankfurt am Main, Hesse, 60438, Germany, <sup>94</sup>Departamento de Botânica, Divisão de Biodiversidade, Instituto Nacional de Pesquisa da Amazônia, Manaus, Amazonas, 69067-375, Brasil, <sup>95</sup>Department of Plant Pathology and Environmental Microbiology, Pennsylvania State University, University Park, State College, PA, 16802, USA, <sup>96</sup>Freshwater Biology Section, Department of Biology, University of Copenhagen, Universitetsparken 4, Copenhagen, 2100, Denmark, <sup>97</sup>Ecosystems and Environment Research Programme, Faculty of Biological and Environmental Sciences, University of Helsinki, Viikinkaari 1, Helsinki, FI-00014, Finland, <sup>98</sup>Department of Zoology, Institute of Ecology and Earth Sciences, University of Tartu, Juhan Liivi 2, Tartu, 50409, Estonia, <sup>99</sup>Department of Plant Sciences, Quaid-i-Azam University, Islamabad, 45320, Pakistan, <sup>100</sup>Department of Agricultural, Food and Forest Sciences, University of Palermo, Viale Delle Scienze, Palermo, 90123, Italy, <sup>101</sup>Department of Medical Physics, College of Science, University of Mosul, Mosul, Iraq, <sup>102</sup>Estonian Museum of Natural History, Lai 29a, Tallinn, 10133, Estonia, <sup>103</sup>Marine ecology, Centre for Advanced Studies of Blanes, CSIC, Carrer Accés Cala Sant Francesc, 14, Blanes, Girona, 17300, Spain, <sup>104</sup>Department of Botany, Jawaharlal Nehru Rajkeeya Mahavidyalaya, Pondicherry University, Port Blair, Andaman and Nicobar Islands, 744101, India, <sup>105</sup>Natural History Museum of Zimbabwe, cnr Park Road & Leopold Takawira Avenue Centenary Park, Bulawayo, Zimbabwe, <sup>106</sup>Department of Geographical and Historical Studies, University of Eastern Finland, Yliopistokatu 7, Joensuu, 80101, Finland, <sup>107</sup>Department of Arctic and Marine Biology, UiT The Arctic University of Norway, Framstredet 39, Tromsø, 9019, Norway, <sup>108</sup>Department of Forestry and Wildlife Management, Science agronomique Kigani-Dada de Kpéssou, Private Agricultural University, Parakou, Benin, <sup>109</sup>Royal Botanic Gardens Victoria, Melbourne, VIC, 3004, Australia, <sup>110</sup>School of BioSciences, University of Melbourne, Parkville, VIC, 3010, Australia, <sup>111</sup>South African Environmental Observation Network, Arid Lands Node, South African Environmental Observation Network, 97 Memorial Road, Kimberley, 8036, South Africa, <sup>112</sup>Department of Biological Sciences, University of Cape Town, Cape Town, 7701, South Africa, <sup>113</sup>UFR Biosciences, Université Félix Houphouët-Boigny, Abidjan, BPV34, Cote d'Ivoire, <sup>114</sup>Grupo BioMicro, Escuela de Microbiología, Universidad de Antioquia UdeA, Calle 70 No. 52-2, Medellín, 50010, Colombia, <sup>115</sup>Department of Invasion Ecology, Institute of Botany of the Czech Academy of Sciences, Zámek 1, Průhonice, 25243, Czechia, <sup>116</sup>Center for Yunnan Plateau Biological Resources, Protection and Utilization & Yunnan International Joint Laboratory of Fungal Sustainable Utilization in South and Southeast Asia, Biology and Food Engineering, Qujing Normal University, Qujing, Yunnan, 655011, China, <sup>117</sup>Department of Microbiology, Anatomy, Physiology and Pharmacology, La Trobe University, Bundoora, VIC, 3086, Australia, <sup>118</sup>Tropical Mycology and Plant-Soil fungi Interactions (MyTIPS), Faculty of Agronomy, University of Parakou, Parakou, Benin, <sup>119</sup>Applied Science Department, William & Mary, 540 Landrum Drive, Williamsburg, VA, 23185, USA, <sup>120</sup>Mycology and Microbiology Center, Institute of Technology, University of Tartu, Nooruse 1, Tartu, 50411, Estonia

Corresponding authors: Vladimir Mikryukov, Leho Tedersoo; ✉e-mail: vladimir.mikryukov@ut.ee, leho.tedersoo@ut.ee

## Abstract

Soil eukaryotes, including fungi, protists, plants, and animals, are central to biosphere functioning and resilience. The Global Standardised Soil Eukaryome Dataset (GloSED) is the first dataset encompassing the entire spectrum of soil eukaryotes, covering 4,063 sampling sites in 121 countries on all continents, revealing nearly one million operational taxonomic units. All samples were collected and analysed using a standardised protocol minimizing technical biases. Long-read sequencing of full-length ITS and 18S-V9 regions provide broad taxonomic coverage and high-resolution identification supported by specialist curation of “dark taxa”. A rigorous bioinformatic processing ensures against homopolymer errors, PCR-mediated chimeras, and index switching providing high data quality. The dataset is supported by raw sequences and an open-source containerised workflow for reproducible analyses. The samples are accompanied by land-cover description and directly measured soil pH,  $\delta^{13}\text{C}$ ,  $\delta^{15}\text{N}$ , as well as P, K, Ca, Mg, and total C and N contents. GloSED is the first database that enables ecological and biogeographic studies of entire soil eukaryotic communities from local to global scales.

## Background & Summary

Soils harbor an extraordinary diversity of eukaryotic life that underpins the functioning of terrestrial ecosystems. Soil-dwelling eukaryotes, including plants, fungi, animals, and diverse groups of protists, are important drivers of nutrient cycling, carbon storage, and ecosystem functioning, stability, and resilience<sup>1</sup>. Over the past decade, global datasets based on environmental DNA sequencing that mapped the diversity and distribution of soil eukaryotes have focused only on individual functional groups or taxa at kingdom, phylum or class levels: fungi in general<sup>2–5</sup> and mycorrhizal fungi in particular<sup>6</sup>, dominant protist taxa<sup>7,8</sup>, nematodes<sup>9</sup>, earthworms<sup>10</sup>, and springtails<sup>11</sup>, among others.

Despite remarkable progress in mapping specific groups, no global dataset has yet captured the full spectrum of eukaryotes in soil samples. Soil organisms, however, do not operate in isolation; instead, they form complex interaction networks, including predation, parasitism, competition, and symbiosis<sup>1,12</sup>. The need to study these groups together within an integrated framework has long been recognised<sup>12–14</sup>. Yet, cross-group comparisons demand standardized protocols, the use of universal metabarcoding primers and curation of taxonomic annotations by specialists in many taxa. Due to these challenges, a globally standardized dataset spanning multiple eukaryotic

kingdoms has not yet been produced, limiting the ability to conduct integrated analyses across taxonomic and functional groups.

Some existing databases, however, offer the potential for expansion. The Global Soil Mycobiome Consortium (GSMc), established in 2014, initially produced a comprehensive dataset of global soil fungal diversity<sup>5</sup>. GSMc employs long-read sequencing of universal eukaryotic primers – full-length internal transcribed spacer (ITS) and partial 18S rRNA gene (SSU) regions – enabling a taxonomic expansion beyond fungi. Here, we leveraged this capability to encompass a broader range of soil eukaryotes, including fungi, protists, animals, and plants. To achieve this, we (i) reprocessed all sequences using an improved bioinformatic workflow to enhance data accuracy; (ii) increased taxonomic resolution by applying the curated EUKARYOME reference database<sup>15</sup>; and (iii) incorporated 1059 new sampling sites, extending the dataset's geographical coverage from 108 to 121 countries.

Here we present a globally representative dataset of soil eukaryotic diversity – GloSED (Global Standardised Soil Eukaryome Dataset). GloSED inherits from its predecessor, the GSMc database, the key advantage of methodological consistency across geographic regions, achieved through standardised sampling and analytical protocols. The dataset is accompanied by soil chemical and habitat metadata and is supported by raw long-read sequence data and a fully automated open-source bioinformatics pipeline that runs in a standardised, portable software environment, ensuring transparency, reproducibility, and flexible user customization.

## Methods

### *Sampling and sample pre-processing*

Soil samples were collected following a standardised approach<sup>2</sup>. Each sampling plot comprised a 50 × 50 m square or a 56 m diameter circular area (2500 m<sup>2</sup>). Within each plot, 40 soil cores (5 cm diameter, 5 cm depth) were collected, with cores taken in pairs on opposite sides of randomly selected trees 1.5 m from the tree trunk (in forests) or at random locations (in non-forested ecosystems). Sampling points were positioned at least 8 m apart to ensure spatial independence while providing comprehensive plot coverage.

Individual soil cores were pooled by combining equal volumes (approximately 25% of the total volume per core), thoroughly mixed, and air-dried within 24 hours of collection. Dried samples were manually homogenised through vigorous rubbing in sealed plastic bags. After homogenization, approximately 30–50 g of the finest material was retained for further analyses. Samples were either shipped to the University of Tartu, Estonia, with silica gel for centralised

processing or processed immediately in contributor laboratories following standardised protocols described below.

All sampling was conducted under appropriate national permits and followed local regulations for soil collection. The DNA extracts and soil samples are stored in the Collections of DNA and environmental samples (TUE) of the Natural History Museum at the University of Tartu.

### *Soil chemical analyses*

Soil pH was measured potentiometrically in 1M KCl extract at a 1:2.5 soil:solution ratio. Available phosphorus and potassium were measured in 1M ammonium lactate extracts by flow injection analysis using a Tecator autoanalyser (method ASTN 9/84). Exchangeable magnesium and calcium were determined in 1M ammonium acetate extracts by flow injection analysis (method ASTN 90/92). At least 0.1 g of ball-milled soil was analysed for total carbon,  $^{13}\text{C}$ , total nitrogen and  $^{15}\text{N}$  content using an elemental analyzer coupled with an isotope ratio mass spectrometer in 2–6 replicates.  $\delta^{13}\text{C}$  and  $\delta^{15}\text{N}$  were obtained following international standards<sup>16</sup>.

### *Molecular analyses*

DNA extraction was performed from 2 g of homogenised dry soil using the PowerMax Soil DNA Isolation kit (Qiagen, Carlsbad, CA, United States) following the manufacturer's protocols. Extracted DNA was further purified using the FavorPrep™ Genomic DNA Clean-Up kit (Favorgen, Vienna, Austria) to remove inhibiting compounds and improve amplification success.

Polymerase chain reaction (PCR) was used to amplify the full internal transcribed spacer (ITS) region and 18S-V9 variable region using universal eukaryotic primers ITS9mun (5'-GTACACACCGCCCGTCG-3') and ITS4ngsUni (5'-CGCCTSCSCTTANTDATATGC-3')<sup>17,18</sup>. This primer combination amplifies nearly all known eukaryotes with minimal taxonomic bias<sup>19</sup>, excluding only Microsporidia due to primer mismatches. Each primer pair was labeled with identical 12-base indices selected from 115 combinations to minimize index-switching artefacts (Hamming distance between each pair of indices was no less than 3).

PCR reactions contained 5  $\mu\text{l}$  of 5x HOT FIREPol Blend Master Mix (Solis Biodyne, Tartu, Estonia), 0.5  $\mu\text{l}$  each of forward and reverse primers (20  $\mu\text{M}$ ), 1  $\mu\text{l}$  DNA extract, and 18  $\mu\text{l}$  ddH<sub>2</sub>O in 25  $\mu\text{l}$  total volume. Thermal cycling comprised initial denaturation at 95 °C for 15 min, followed by 25-30 cycles of denaturation (30 s at 95 °C), annealing (30 s at 57 °C), and elongation (1 min at 72 °C), with final extension at 72 °C for 10 min. PCR products were checked on 1% agarose gels for 600-800 bp amplicons. Samples failing initial amplification were re-amplified using 28-38 cycles.

Library preparation and sequencing were performed at the Norwegian Sequencing Centre, University of Oslo, Norway. PacBio SMRTbell libraries were prepared following the manufacturer's protocols and sequenced on Sequel II instruments using Sequel II Binding kit 2.1, sequencing chemistry 2.0, with 15-hour movie times and 20-minute pre-extension periods. Samples producing <2000 reads after the first sequencing attempt were re-amplified and re-sequenced (this applied to 929 samples). For these samples, reads from repeated runs were combined in the final dataset (in raw sequencing data, reads are provided in independent files for each sequencing attempt), and 97.8% reached the minimum target depth. Sequencing was distributed across 122 SMRT cells.

### ***Bioinformatics pipeline***

GloSED was analysed using the fully automated bioinformatics pipeline NextITS v.1.0.0<sup>20</sup> (DOI:10.5281/zenodo.15074882) implemented with the workflow manager Nextflow v.25.04.6<sup>21</sup> and run in standardised, portable software environments (Docker and Singularity containers)<sup>22,23</sup> to ensure reproducibility. The open-source NextITS pipeline is freely available as a command-line workflow at <https://Next-ITS.github.io/> and is also distributed through the cross-platform PipeCraft2 application<sup>24</sup> (<https://pipecraft2-manual.readthedocs.io/en/latest/>), which provides a graphical user interface for running the pipeline.

Raw reads were processed through circular consensus sequence (CCS) generation using SMRT Tools (Pacific Biosciences) with minimum pass requirements of 3 and an accuracy threshold of 0.99. Demultiplexing was performed using LIMA v.2.12.0 (Pacific Biosciences) with `--min-score 93` settings for precise index identification. Quality filtering removed sequences with >4 ambiguous nucleotides, >0.01% expected errors<sup>25</sup>, and homopolymer repeats longer than 25 nucleotides. In addition, reads lacking both primer sites in the correct orientation were discarded after primer trimming with cutadapt v.5.0<sup>26</sup>. Full-length ITS regions were then retrieved by ITSx v.1.1.3<sup>27</sup> targeting all eukaryotes using the updated hidden Markov model (HMM) profile database (provided by R. Henrik Nilsson, University of Gothenburg, Sweden). For sequence processing, we used SeqKit2 v.2.9.0<sup>28</sup>. When selecting representative ITS sequences after ITSx-based extraction, preference was given to the sequence with the highest average Phred score. Chimeras were detected in a two-step scheme with VSEARCH v.2.29.4<sup>29</sup>: an initial *de novo* detection using the UCHIME algorithm<sup>29</sup> and a maximum chimera score of 0.6<sup>30</sup> followed by reference-based verification against the EUKARYOME v.1.9.4<sup>15</sup> database, with any sequence flagged in either step being excluded. Within each sample, homopolymer correction was performed using the algorithm implemented in NextITS. Within each sequencing run, index-switch (tag-jump) artefacts were removed using the UNCROSS2 algorithm<sup>31</sup> with the parameter  $f = 0.01$ . Prior to clustering, ITS sequences shorter than 250 bp or containing more than 0.6 expected errors per 100 bp were discarded, and the remaining reads were denoised using the UNOISE3 algorithm<sup>32</sup> with parameters

alpha = 6 and minsize = 1, which retains singleton sequences. Surviving denoised reads were then clustered at 98% pairwise similarity using VSEARCH, and these clusters were used as operational taxonomic units<sup>33</sup> (OTUs) in downstream analyses, approximating species-level groupings and collapsing individual, potentially intragenomic sequence variants<sup>34,35</sup>. We thus retain rare variants at the denoising step and summarise diversity at the 98%-similarity OTU level, rather than treating each denoised exact sequence variant as a separate analytical unit, because long amplicons contain many singleton and other low-abundance variants that may represent genuine rare taxa, and removing them at the denoising step could lead to an underestimation of rare diversity<sup>34,36</sup>. The 98% similarity threshold represents a pragmatic compromise across divergent eukaryotic lineages and may affect richness estimates in groups with different rates of ITS evolution, but implementing lineage-specific clustering thresholds at this scale is currently not feasible.

Representative ITS sequences of each OTU, as well as the corresponding small (SSU) and large (LSU) subunit fragments of the same representative, were queried with BLASTn v.2.16.0+<sup>37</sup> against the EUKARYOME database, retaining the ten best hits. Extensive manual curation was performed using taxon-specific E-value and sequence similarity thresholds<sup>5</sup>. Additionally, long-read chimeras were detected by comparing region-specific taxonomic assignments across the SSU, ITS, and LSU segments of each read. Reads were classified as chimeric when these segments yielded incongruent higher-level placements, as determined by manual inspection of the top BLAST hits, using phylum-level disagreement for Ascomycota and Basidiomycota and order-level disagreement for other taxa. Disagreement restricted to lower taxonomic ranks was not treated as sufficient evidence of chimerism. Non-target OTUs (e.g., sequences of archaeal or bacterial origin) were removed. Samples with <500 total reads or identified as potentially contaminated (i.e., had >30% of reads corresponding to molds or shared OTUs with negative or positive controls) were excluded from the final dataset. To reduce the amount of unidentified OTUs and improve the precision of chimera detection, we also amplified and sequenced an ultra-long rRNA gene fragment spanning the SSU V3 through ITS through LSU D8 from >900 soil samples using the primers EUK575F (5'-TASCYGYGGTAAYWCCAGC-3') and 21R (5'-AGAGACGAGGCATTTGGCTAC-3')<sup>38,39</sup> on PacBio Sequel II and Revio platforms.

To complement EUKARYOME-based taxonomic annotations of fungal OTUs, we additionally performed UNITE Species Hypotheses (SH) matching<sup>40-42</sup> to assign persistent DOI-registered identifiers to representative ITS sequences. The analysis was performed using the SH-matching analysis tool v.2.0.3, which places query ITS sequences into existing SHs in the UNITE database v.10.0<sup>43</sup> or assigns to new SHs with preliminary codes. The stable identifiers enable unambiguous cross-study referencing of “dark taxa”<sup>44</sup>, remain valid as taxonomy changes, and facilitate reporting and reuse of DNA-derived occurrences. Because SHs are integrated into the GBIF

(<https://www.gbif.org/>) taxonomic backbone, SH-based occurrences from this study can be directly compared, aggregated, and cross-linked with GBIF records and other datasets.

## Data Records

The GloSED dataset is available at Zenodo<sup>45</sup> (versioned release archive DOI: 10.5281/zenodo.17827890). The Zenodo deposit is organised as a set of files, including: (1) detailed sample metadata with environmental variables ('GloSED\_\_Sample\_metadata.xlsx'), (2) quality-filtered OTU sequences in FASTA format ('GloSED\_\_OTU\_sequences.fasta.gz'), (3) sample-by-OTU abundance matrices ('GloSED\_\_OTU\_table.tsv.zip' and 'GloSED\_\_OTU\_table.parquet'), and (4) curated taxonomic annotations ('GloSED\_\_Taxonomy.tsv.zip' and 'GloSED\_\_Taxonomy.parquet').

'GloSED\_\_Sample\_metadata.xlsx' is the main metadata table and contains one row per sample. The accompanying legend sheet defines all column names and measurement units. The key linking fields are 'SampleID' and 'TUE code', where the latter is the accession of the physical soil sample in the University of Tartu collection. The metadata describe when and where each sample was collected, including geographic coordinates, elevation, locality, administrative unit, land-cover type, and co-occurring plant taxa. The table also includes measured soil properties, including pH, total soil carbon and nitrogen, C:N ratio, stable isotope abundances, as well as nutrient concentrations for phosphorus, potassium, calcium, and magnesium.

OTU abundance data are provided as 'GloSED\_\_OTU\_table.tsv.zip' and 'GloSED\_\_OTU\_table.parquet'. The Parquet file stores the abundance information in long format with the fields 'OTU', 'SampleID', and 'Abundance', whereas the tab-delimited file contains sample-by-OTU abundance matrix in wide format. Taxonomic annotation is provided in 'GloSED\_\_Taxonomy.tsv.zip' and 'GloSED\_\_Taxonomy.parquet'. These files contain the OTU identifier, the representative sequence, accession number of the best EUKARYOME match (column 'AccID') and its alignment statistics (sequence identity percentage, coverage, E-value, and BLAST bitscore), curated taxonomic ranks from 'Kingdom' to 'Species', and the UNITE fields 'SH\_30' to 'SH\_05', which record persistent fungal species-hypothesis identifiers at progressively finer similarity thresholds (from 3% to 0.5%). Files in Apache Parquet<sup>46</sup> format provide a programming-language-independent columnar file structure that allows efficient storage and high-performance analytical operations, whereas TSV files offer broader compatibility. Representative OTU sequences are supplied separately in 'GloSED\_\_OTU\_sequences.fasta.gz', with FASTA headers matching the OTU identifiers (40-character hexadecimal strings based on the SHA1 hash of the primer-trimmed sequence prior to ITSx extraction).

For common downstream workflows, the same core data are additionally distributed as a biological observation matrix (BIOM) v.2.1 format file<sup>47</sup> for QIIME2<sup>48</sup> integration (file

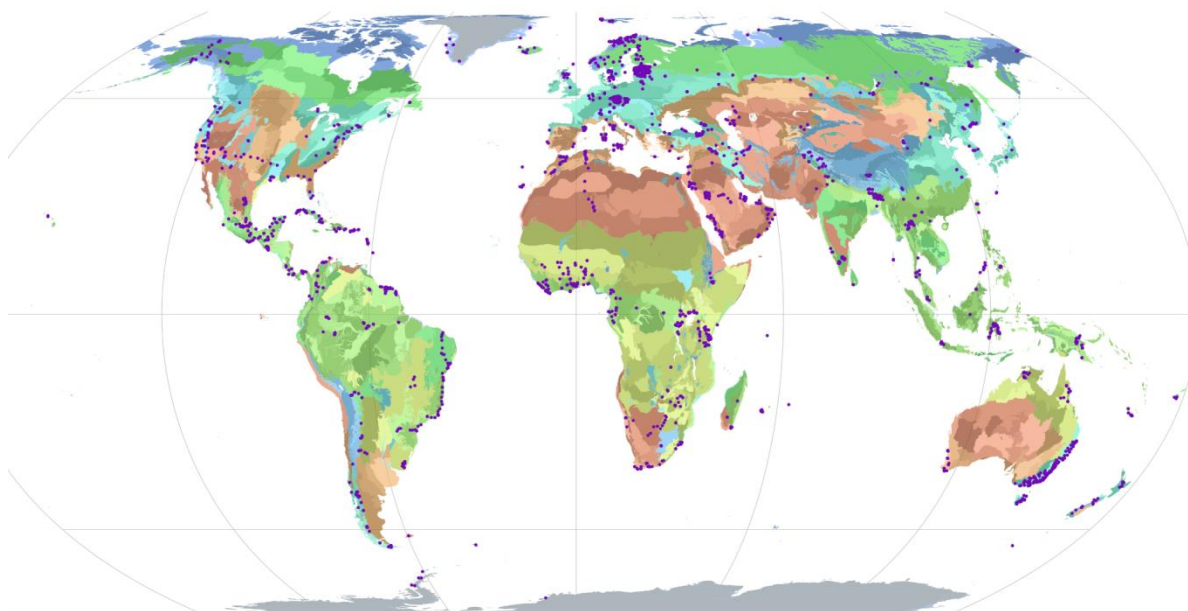
'GloSED\_\_BIOM.biom') and as a *phyloseq*<sup>49</sup> object for R-based analyses ('GloSED\_\_phyloseq.RData'). 'Contributors.xlsx' provides the contributor list associated with the release. Raw demultiplexed FASTQ files are available from the European Nucleotide Archive (ENA) under project PRJEB10381182<sup>50</sup>.

### Data Overview

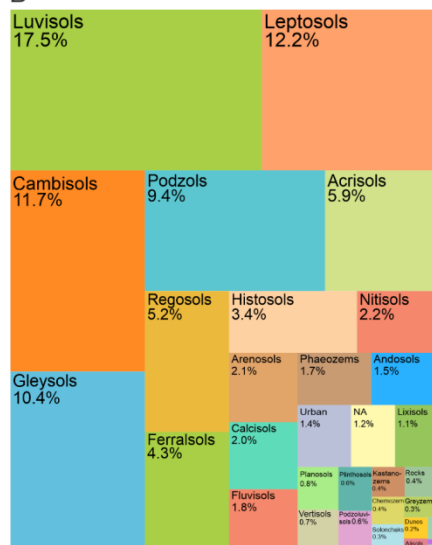
The GloSED is a structured dataset including 4,147 samples from 4,063 sampling sites worldwide (Fig. 1A), analysed using standardised field and laboratory protocols. Each site record contains geographic coordinates, sampling date, land-cover type assigned by the field collector, and plot-level information on the dominant plants, mostly for woody land cover types. For nearly 95% of sites, directly measured soil properties are recorded, including pH, total carbon and total nitrogen contents ( $\text{g kg}^{-1}$ ),  $\delta^{15}\text{N}$  (‰),  $\delta^{13}\text{C}$  (‰), as well as available phosphorus and potassium, and exchangeable magnesium and calcium contents (each  $\text{mg kg}^{-1}$ ).

All samples are associated with read counts for each of 988,824 operational taxonomic units (OTUs) of eukaryotic organisms. Each OTU is taxonomically annotated using the EUKARYOME reference database, assigned a species hypothesis (SH) identifier using the UNITE reference database, and linked to quality-curated representative sequences of full-length ITS and 18S-V9 regions. Median observed OTU richness and effective number of OTUs per sample were 829 (Q1-Q3: 526-1213) and 184 (Q1-Q3: 106-280), respectively.

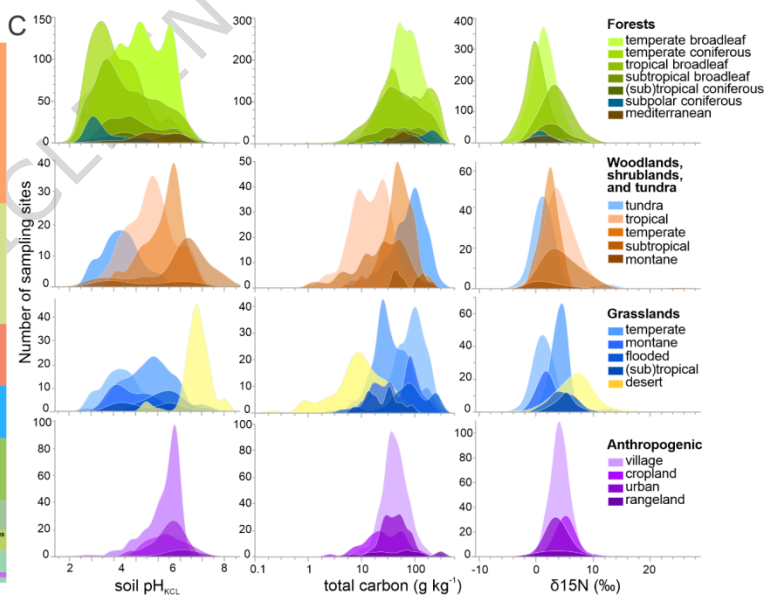
A



B



C



**Fig. 1. Examples of environmental gradients represented by GloSED (A)** GloSED sampling sites are distributed across 351 terrestrial ecoregions. Dots mark sampling sites; ecoregions within the same biome share color hues. **(B)** GloSED represents soil types in proportion to their global distribution. **(C)** Soil pH, total carbon, and  $\delta^{15}\text{N}$ , measured directly in GloSED samples, have broad distributions within and across land-cover types.

## Technical validation

### *Methods for technical data validation*

All analyses, and visualization were conducted in R v.4.5.2<sup>51</sup>. For high-performance operations on large datasets, we performed data manipulation and processing using the `data.table` v.1.17.8<sup>52</sup> and `Apache arrow` v.21.0.0<sup>53</sup> packages. Taxonomic resolution and sequencing completeness metric were assessed with the `metagMisc` package v.0.5.0<sup>54</sup>. Observed abundance-based sample coverage was used to estimate sequencing completeness after correcting singleton counts with a modified Good-Turing estimator<sup>55</sup>. The effective number of OTUs was calculated as the exponential of the Shannon diversity index, following Jost<sup>56</sup>.

Spatial data were processed using the `sf` package v.1.0-21<sup>57</sup>. For map visualizations, we used the Equal Earth projection<sup>58</sup>. Sample administrative location was assigned using GADM v.4.1 boundary polygons<sup>59</sup>, biomes and ecozones were determined following Loidi et al.<sup>60</sup>, and ecoregions and soil types were assigned using the Ecoregions 2017 dataset<sup>61</sup> and Harmonized world soil database v.2.0<sup>62</sup>, respectively. We quantified environmental novelty of unsampled locations using a dissimilarity index (DI) following the area-of-applicability framework<sup>63</sup>. The DI was calculated from 14 Z-standardized, importance-weighted bioclimatic and edaphic predictors identified as key drivers of soil fungal communities<sup>4</sup>.

### *Sampling consistency and representativeness*

A consistent sampling design is a fundamental prerequisite for obtaining unbiased biodiversity data<sup>64</sup>. For soil biota in particular, differences in sampling time, effort, area, depth at which soil is collected, and compositing strategy can cause orders-of-magnitude variation in biodiversity estimates<sup>65,66</sup>. Similar effects arise from inconsistencies in analytical protocols<sup>65-67</sup>. GloSED minimizes these potential biases by applying identical sampling and analytical procedures across all sites.

GloSED applies a 5 cm sampling depth to target the most biologically active soil layer where root density, microbial activity, and organic matter content are typically the highest<sup>68</sup>. The relatively large sampling area (2,500 m<sup>2</sup>) – one of the largest commonly used in soil microbial ecology – ensures comprehensive representation of local biota, including soil macrofauna such as annelids and arthropods.

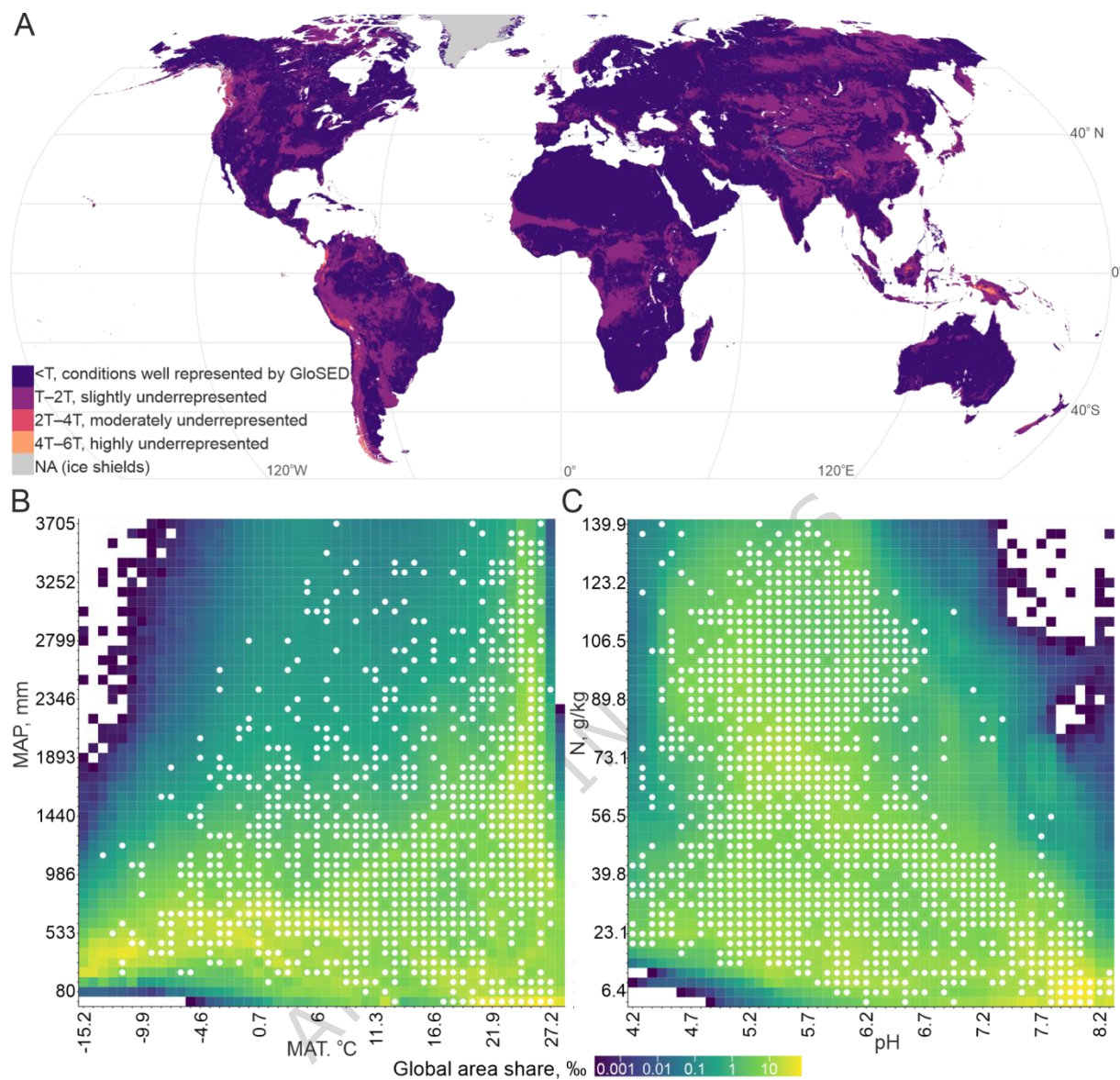
Furthermore, on-site descriptions of land cover by collectors coupled with direct measurements of soil chemical properties, provide data that are more representative of real-world conditions than those derived from modelled or predicted datasets.

### ***Bioinformatic validation***

For the database, more than 73 million HiFi PacBio DNA reads were generated from 122 sequencing runs targeting full-length ITS and 18S-V9 variable regions. During bioinformatic processing, to ensure high data quality, we paid special attention to homopolymer errors, PCR-mediated chimeras, and index switching - technical artefacts that are common in amplicon sequencing but typically undetectable by default pipelines<sup>18,69</sup>. Comparison with reference sequences revealed more than one million chimeric reads, and *de novo* analysis combined with manual curation flagged additional 58,119 chimeras. Index switching rate was within acceptable thresholds (0.02%). All the chimeric reads and false assignments were removed prior to downstream analyses. The 4,147 samples that passed all stages of technical control yielded more than 27 million high-quality sequences with a median number of 5,193 sequences per sample (IQR: 2,766–8,409). Median sample coverage was 0.95 (Q1–Q3: 0.92–0.97), suggesting a high sequencing completeness for most samples.

### ***Geographical and environmental coverage***

Sufficient geographical and environmental coverage is a prerequisite for a global dataset suitable for cross-regional comparisons and predictive modelling<sup>63,70</sup>. The GloSED spans from  $-73.0^{\circ}$  to  $79.6^{\circ}$  latitude, encompassing all terrestrial climate zones (as per Beck et al.<sup>71</sup>), 351 terrestrial ecoregions and all major soil types with broad distributions of soil properties (Fig. 1). Previous global models of soil biodiversity identified regions with high uncertainty in predictions, as these regions differ substantially in environmental conditions from the data observed at sampling locations<sup>4,72</sup>. In a recent sampling campaign, we specifically targeted such regions. For instance, substantial contributions of samples from Australia, India, Saudi Arabia, Algeria, Liberia, Kazakhstan, Suriname, and Uruguay improved the representativeness of hot and temperate deserts as well as tropical and grassland biomes (Fig. 2) – large areas that have shown the greatest variability in soil fungal biodiversity predictions across studies.



**Fig. 2. Examples of GloSED coverage of global environmental space.** (A) Based on 14 climatic, vegetation, and edaphic predictors of total soil fungal diversity<sup>4</sup>, 71.6% of global terrestrial area falls within the GloSED sampling range (dark purple), and 26.3% deviates from this range by less than twofold (purple). Red and orange colors (0.72% and 0.04% of total area) denote the most environmentally novel terrestrial areas relative to the currently sampled space. (B) and (C) show that GloSED captures the majority of the global range of mean annual temperature (MAT) with precipitation (MAP), and topsoil pH with total nitrogen concentration, respectively. Dots mark environmental bins represented by GloSED samples.

### ***Taxonomic coverage and resolution***

Owing to the use of universal eukaryotic markers (full-length ITS and 18S-V9 regions), GloSED is the only pan-domain soil dataset, with a standardized design, currently available. It encompasses 546,665 fungal OTUs (73.6% of reads), 316,021 protistan OTUs (16.8% of reads), 103,021 animal OTUs (3.4% of reads), and 19,186 plant OTUs (5.7% of reads) (Fig. 3A).

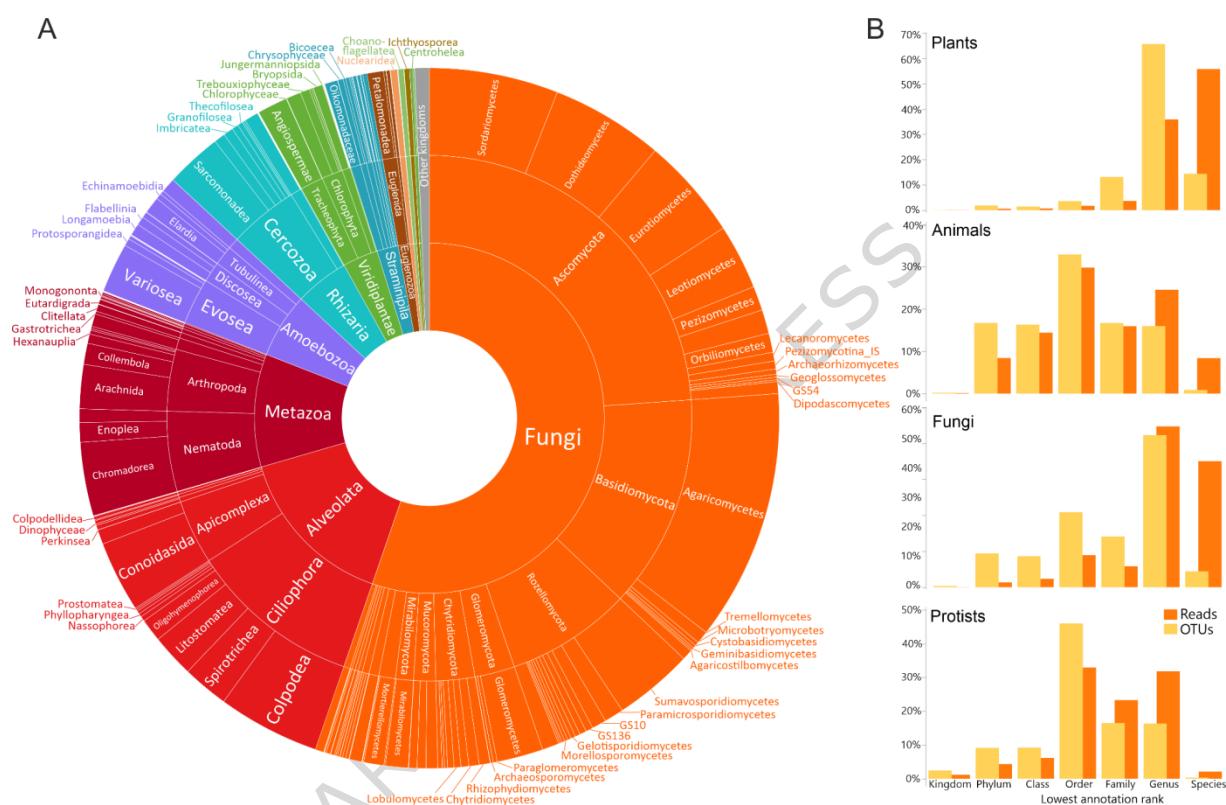
Over 40% of fungal OTUs were taxonomically annotated to the genus level (Fig.3B). As is typical for soil fungal communities, Ascomycota (235,918 OTUs) and Basidiomycota (130,119 OTUs, with 88% belonging to Agaricomycetes) represent the richest phyla. Notably, GloSED taxonomic annotation based on EUKARYOME, curated by leading specialists in fungal and protist taxonomy<sup>73-76</sup>, achieves exceptional resolution for underexplored lineages that may account for a substantial part of the community. For instance, non-Dikarya fungi comprise 180,367 OTUs (33% of total fungal diversity), including 76,113 OTUs from Rozellomycota (also known as Cryptomycota), 24,388 OTUs from Glomeromycota, 23,545 OTUs from Chytridiomycota, 13,319 OTUs from Mucoromycota, and 10,088 OTUs from Mortierellomycota. Owing to a substantially redesigned bioinformatic workflow, the sequence representatives in GloSED differ in both composition and number from those in GSMc.

Among animals, almost half of the data belong to Nematoda (45.6% of metazoan reads; 48,674 OTUs), followed by Arthropoda (36.0% of reads; 37,545 OTUs) and Annelida (10.3% of reads; 3,432 OTUs). Other represented groups include Gastrotricha (2.3% of metazoan reads; 7,097 OTUs), Tardigrada (2.9%; 2,522 OTUs), Rotifera (0.6%; 1,624 OTUs), and Platyhelminthes (1.6%; 1,511 OTUs). Over 16% of animal OTUs are taxonomically annotated to the family level and over 32%, to the order level.

Protistan data span more than 30 phyla, with 74.9% of reads belonging to Alveolata (150,508 OTUs), including almost 40,000 Apicomplexa OTUs - a largely parasitic group associated with major public-health burdens and substantial economic losses<sup>77</sup>. Almost 7% of protist reads belong to each of Amoebozoa (59,179 OTUs) and Rhizaria (49,195 OTUs), 4.7% to green algae (13,120 OTUs), and 3.8% to Straminipila (19,580 OTUs) with 3,308 OTUs belonging to Oomycota – the most economically damaging protist group<sup>78</sup>. Almost 50% of protist OTUs are identified to the order level.

Tracheophyta (vascular plants) cover 73.1% of plant reads (14,030 OTUs, with 98% belonging to angiosperms), and Setaphyta (bryophytes) 26.9% (5,156 OTUs). Over 65% of plant OTUs are annotated to the genus level. This enables the inclusion of plants - key ecosystem engineers – alongside soil microbes, reflecting their reciprocal influence as both drivers and responders within soil biodiversity patterns.

The broad taxonomic coverage of GloSED enables integrative research of soil eukaryotes including primary producers, decomposers, consumers, and parasites. The achieved taxonomic resolution allows functional annotation for a substantial proportion of the GloSED OTUs using databases such as FungalTraits<sup>79</sup> for fungi, Nemaplex<sup>80</sup> for nematodes, FunctionalTraitsAmoebozoa<sup>81</sup> for Amoebozoa and Rhizaria, and TRY for plants<sup>82</sup>.



**Fig. 3. GloSED taxonomic coverage and resolution.** (A) GloSED taxonomic coverage and (B) resolution allows for functional annotation of the majority of soil eukaryotes. Sector angular span denotes the percentage of operational taxonomic units (OTUs) belonging to taxonomic groups. Due to limited space, not all taxa are labelled.

### Cross-dataset comparability

As taxonomy constantly advances, OTUs must be traced across multiple data sources supporting their re-annotation. This is achieved with the UNITE species hypothesis (SH) system – a standardised digital framework for discovering and communicating fungal species, particularly those identified from environmental DNA<sup>42</sup>. In the GloSED dataset, nearly half of all fungal OTUs are assigned to existing SHs at the 97% similarity cut-off, with approximately 20% assigned at the 99.95% threshold. Each SH is associated with a stable DOI-linked reference integrated with the PlutoF and GBIF taxonomic backbones.

## Usage Notes

Raw sequencing data can be re-analysed using the NextITS workflow. The source code of this bioinformatics pipeline and its documentation are available under MIT license at <https://github.com/vmikk/NextITS> and <https://Next-ITS.github.io/>, respectively. For reproducible execution, containerised environments are hosted at Docker Hub (<https://hub.docker.com/r/vmikk/nextits>) and Singularity library (<https://cloud.sylabs.io/library/vmiks/nextits/nextits>). Processing the full GloSED dataset on the University of Tartu high-performance computing (HPC) cluster using AMD EPYC 7702 processors required approximately 13,300 CPU hours, with wall-clock runtime depending on the execution profile and degree of parallelisation. ITSx represented the main bottleneck in the pipeline, and the BLAST searches required an additional approximately 20,600 CPU hours. To support data reuse, we followed the guidelines of Hug et al. (2025)<sup>83</sup>.

## Data Availability

The dataset and sample metadata are available on Zenodo<sup>45</sup> (<https://zenodo.org/records/17827890>), and the raw sequencing data have been deposited in the European Nucleotide Archive (ENA) under project accession PRJEB103811 (sample accession numbers ERS27941879 - ERS27946063; sequence accession numbers ERR15957609 - ERR15964175)<sup>50</sup>.

## Code Availability

Analysis scripts for manuscript figures and supporting analyses are deposited at GitHub (<https://github.com/Mycology-Microbiology-Center/GloSED>).

## Acknowledgments

We thank all researchers, taxonomic specialists, and volunteers who provided help and local expertise essential for this global sampling effort. We are grateful to Jane Lees for her dedicated assistance and for facilitating many aspects of this project, and to the anonymous reviewers for their constructive comments. We acknowledge the computational resources provided by the high performance computing center of the University of Tartu and the sequencing facilities at the Norwegian Sequencing Centre, University of Oslo. We are grateful to CapeNature (Western Cape Nature Conservation Board, South Africa) and the relevant permitting authorities for providing the necessary research and collection permits.

## Author Contributions

L.T. conceived and designed the study, developed the data collection protocol and sampling strategy, led the global soil-sampling collaboration, curated reference databases and verified taxonomic annotations, administered the project, secured funding, and oversaw budget allocation. L.T., O.D, and V.M. performed data cleaning and harmonization. V.M. developed the bioinformatics pipeline and processed the sequencing data. O.D. and V.M. conducted statistical and geospatial analyses, generated figures, maps, and graphical summaries. R.P., P.P., V.P., N.H., and E.O. extracted DNA and prepared libraries for sequencing. V.M. prepared the public data and code repositories. Ke.Ab. developed species hypothesis matching pipeline and managed the dataset hosted on the PlutoF system. S.A. and Ke.Ab. provided metabarcoding domain expertise. The remaining co-authors primarily contributed to field sampling and metadata collection. V.M. and O.D. wrote the original draft of the manuscript. All authors reviewed and edited the manuscript.

## Competing Interests

The authors declare no competing financial or non-financial interests.

## Funding

This work was supported by the Estonian Research Council (projects PRG632, MOBERC116, and PRG1789), the Estonian Ministry of Education and Research (Center of Excellence AgroCropFuture “Agroecology and new crops in future climates”, TK200), the European Research Council (ERC) grant 101200758 (PhylFun ERC-2024-ADG), and Ongoing Research Funding program (ORF-2026-26) by King Saud University (Riyadh, Saudi Arabia). AA acknowledges financial support from the Swedish Research Council (2024-04303), the Swedish Foundation for Strategic Environmental Research MISTRA (project BioPath) and RBG Kew Development.

## References

1. Anthony, M. A., Bender, S. F. & van der Heijden, M. G. A. Enumerating soil biodiversity. *Proc. Natl. Acad. Sci.* **120**, e2304663120 (2023).
2. Tedersoo, L. *et al.* Global diversity and geography of soil fungi. *Science* **346**, 1078 (2014).
3. Větrovský, T. *et al.* GlobalFungi, a global database of fungal occurrences from high-throughput-sequencing metabarcoding studies. *Sci. Data* **7**, 228 (2020).
4. Mikryukov, V. *et al.* Connecting the multiple dimensions of global soil fungal diversity. *Sci. Adv.* **9**, eadj8016 (2023).
5. Tedersoo, L. *et al.* The Global Soil Mycobiome consortium dataset for boosting fungal diversity research. *Fungal Divers.* **111**, 573–588 (2021).

6. Větrovský, T. *et al.* GLOBALAMFUNGI : a global database of arbuscular mycorrhizal fungal occurrences from high-throughput sequencing metabarcoding studies. *New Phytol.* **240**, 2151–2163 (2023).
7. Bates, S. T. *et al.* Global biogeography of highly diverse protistan communities in soil. *ISME J.* **7**, 652–659 (2013).
8. Oliverio, A. M. *et al.* The global-scale distributions of soil protists and their contributions to belowground systems. *Sci. Adv.* **6**, eaax8787 (2020).
9. Van Den Hoogen, J. *et al.* A global database of soil nematode abundance and functional group composition. *Sci. Data* **7**, 103 (2020).
10. Phillips, H. R. P. *et al.* Global distribution of earthworm diversity. *Science* **366**, 480–485 (2019).
11. Potapov, A. M. *et al.* Globally invariant metabolism but density-diversity mismatch in springtails. *Nat. Commun.* **14**, 674 (2023).
12. Bardgett, R. D. & van der Putten, W. H. Belowground biodiversity and ecosystem functioning. *Nature* **515**, 505–511 (2014).
13. Bahram, M. *et al.* Structure and function of the global topsoil microbiome. *Nature* **560**, 233–237 (2018).
14. Netherway, T., Bengtsson, J., Krab, E. J. & Bahram, M. Biotic interactions with mycorrhizal systems as extended nutrient acquisition strategies shaping forest soil communities and functions. *Basic Appl. Ecol.* **50**, 25–42 (2021).
15. Tedersoo, L. *et al.* EUKARYOME: the rRNA gene reference database for identification of all eukaryotes. *Database J. Biol. Databases Curation* **2024**, baae043 (2024).
16. Werner, R. A. & Brand, W. A. Referencing strategies and techniques in stable isotope ratio analysis. *Rapid Commun. Mass Spectrom.* **15**, 501–519 (2001).
17. Tedersoo, L. & Lindahl, B. Fungal identification biases in microbiome projects. *Environ. Microbiol. Rep.* **8**, 774–779 (2016).
18. Tedersoo, L. & Anslan, S. Towards PacBio-based pan-eukaryote metabarcoding using full-length ITS sequences. *Environ. Microbiol. Rep.* **11**, 659–668 (2019).
19. Tedersoo, L., Tooming-Klunderud, A. & Anslan, S. PacBio metabarcoding of Fungi and other eukaryotes: errors, biases and perspectives. *New Phytol.* **217**, 1370–1385 (2018).
20. Mikryukov, V., Anslan, S. & Tedersoo, L. NextITS: a pipeline for metabarcoding eukaryotes with full-length ITS sequenced with PacBio. <https://next-ITS.github.io/> (2025).
21. Di Tommaso, P. *et al.* Nextflow enables reproducible computational workflows. *Nat. Biotechnol.* **35**, 316–319 (2017).
22. Kurtzer, G. M., Sochat, V. & Bauer, M. W. Singularity: Scientific containers for mobility of compute. *PLOS ONE* **12**, e0177459 (2017).
23. Merkel, D. Docker: lightweight Linux containers for consistent development and deployment. *Linux J* **2014**, 2:2 (2014).
24. Anslan, S., Bahram, M., Hiiesalu, I. & Tedersoo, L. PipeCraft: Flexible open-source toolkit for bioinformatics analysis of custom high-throughput amplicon sequencing data. *Mol. Ecol. Resour.* **17**, e234–e240 (2017).

25. Edgar, R. C. & Flyvbjerg, H. Error filtering, pair assembly and error correction for next-generation sequencing reads. *Bioinformatics* **31**, 3476–3482 (2015).
26. Martin, M. Cutadapt removes adapter sequences from high-throughput sequencing reads. *EMBnet.journal* **17**, 10–12 (2011).
27. Bengtsson-Palme, J. *et al.* Improved software detection and extraction of ITS1 and ITS2 from ribosomal ITS sequences of fungi and other eukaryotes for analysis of environmental sequencing data. *Methods Ecol. Evol.* **4**, 914–919 (2013).
28. Shen, W., Sipos, B. & Zhao, L. SeqKit2: A Swiss army knife for sequence and alignment processing. *iMeta* **3**, e191 (2024).
29. Rognes, T., Flouri, T., Nichols, B., Quince, C. & Mahé, F. VSEARCH: a versatile open source tool for metagenomics. *PeerJ* **4**, e2584 (2016).
30. Nilsson, R. H. *et al.* A Comprehensive, Automatically Updated Fungal ITS Sequence Dataset for Reference-Based Chimera Control in Environmental Sequencing Efforts. *Microbes Environ.* **30**, 145–150 (2015).
31. Edgar, R. C. UNCROSS2: identification of cross-talk in 16S rRNA OTU tables. 400762 Preprint at <https://doi.org/10.1101/400762> (2018).
32. Edgar, R. C. UNOISE2: improved error-correction for Illumina 16S and ITS amplicon sequencing. 081257 Preprint at <https://doi.org/10.1101/081257> (2016).
33. Blaxter, M. *et al.* Defining operational taxonomic units using DNA barcode data. *Philos. Trans. R. Soc. B Biol. Sci.* **360**, 1935–1943 (2005).
34. Tedersoo, L. *et al.* Best practices in metabarcoding of fungi: From experimental design to results. *Mol. Ecol.* **31**, 2769–2795 (2022).
35. Kausrud, H. ITS alchemy: On the use of ITS as a DNA marker in fungal ecology. *Fungal Ecol.* **65**, 101274 (2023).
36. Cazabonne, J., Walker, A. K., Lesven, J. & Haelewaters, D. Singleton-based species names and fungal rarity: Does the number really matter? *IMA Fungus* **15**, 7 (2024).
37. Camacho, C. *et al.* BLAST+: architecture and applications. *BMC Bioinformatics* **10**, 421 (2009).
38. Schwelm, A., Berney, C., Dixelius, C., Bass, D. & Neuhauser, S. The Large Subunit rDNA Sequence of *Plasmodiophora brassicae* Does not Contain Intra-species Polymorphism. *Protist* **167**, 544–554 (2016).
39. Aslani, F. *et al.* Land use intensification homogenizes soil protist communities and alters their diversity across Europe. *Soil Biol. Biochem.* **195**, 109459 (2024).
40. Kõljalg, U. *et al.* The Taxon Hypothesis Paradigm—On the Unambiguous Detection and Communication of Taxa. *Microorganisms* **8**, 1910 (2020).
41. Kõljalg, U. *et al.* Towards a unified paradigm for sequence-based identification of fungi. *Mol. Ecol.* **22**, 5271–5277 (2013).
42. Abarenkov, K., Kõljalg, U. & Nilsson, R. H. UNITE Species Hypotheses Matching Analysis. in *Biodiversity Information Science and Standards* vol. 6 e93856 (Pensoft Publishers, 2022).

43. Abarenkov, K. *et al.* The UNITE database for molecular identification and taxonomic communication of fungi and other eukaryotes: sequences, taxa and classifications reconsidered. *Nucleic Acids Res.* **52**, D791–D797 (2024).
44. Nilsson, R. H. *et al.* How, not if, is the question mycologists should be asking about DNA-based typification. *MycKeys* **96**, 143–157 (2023).
45. Mikryukov, V., GSMc consortium, Tedersoo, L. GloSED - Global standardised soil eukaryome dataset. *Zenodo* <https://doi.org/10.5281/zenodo.17827890> (2026).
46. Apache Parquet development team. *Apache Parquet: an open source, column-oriented data file format release 1.16.0*. <https://parquet.apache.org/> (2025).
47. McDonald, D. *et al.* The Biological Observation Matrix (BIOM) format or: how I learned to stop worrying and love the ome-ome. *GigaScience* **1**, 2047-217X-1–7 (2012).
48. Bolyen, E. *et al.* Reproducible, interactive, scalable and extensible microbiome data science using QIIME 2. *Nat. Biotechnol.* **37**, 852–857 (2019).
49. McMurdie, P. J. & Holmes, S. phyloseq: An R Package for Reproducible Interactive Analysis and Graphics of Microbiome Census Data. *PLOS ONE* **8**, e61217 (2013).
50. *ENA European Nucleotide Archive*. <https://identifiers.org/ena.embl:PRJEB103811> (2026).
51. R Core Team. *R A Language and Environment for Statistical Computing version 4.5.2*. <https://www.scirp.org/reference/referencespapers?referenceid=3967248> (2025).
52. Barrett T., Dowle M., Srinivasan A., Gorecki J., Chirico M., Hocking T., Schwendinger B., Krylov I., Stetsenko P., Short T., Lianoglou S., Antonyan E., Bonsch M., Parsonage H., Ritchie S., Ren K., Tan X., Saporta R., Seiskari O., Dong X., Lang M., Iwasaki W., Wenchel S., Broman K., Schmidt T., Arenburg D., Smith E., Cocquemas F., Gomez M., Chataignon P., Blaser N., Selivanov D., Riabushenko A., Lee C., Groves D., Possenriede D., Parages F., Toth D., Yaramaz-David M., Perumal A., Sams J., Morgan M., Quinn M., Storey R., Saraswat M., Jacob M., Schubmehl M., Vaughan D., Silvestri L., Hester J., Damico A., Freundt S., Simons D., Sales de Andrade E., Miller C., Meldgaard J.P., Tlapak V., Ushey K., Eddelbuettel D., Fischetti T., Shilon O., Khotilovich V., Wickham H., Becker B., Haynes K., Kamgang B.C., Delmarcell O., O'Brien J., de Mezquita D., Czekanski M., Shemetov D., Jha N., Wu J., Giné-Vázquez I., Chetia A., Amoakohene D., Feliz A., Young M., Seeto M., Grosjean P., Runge V., Wia C., Maigné E., Rocher V., Lulla V., Sluga A., Evans B. *data.table: Extension of 'data.frame'* <https://cran.r-project.org/web/packages/data.table/index.html> (2025).
53. Richardson, N., Cook I., Crane N., Dunnington D., François R., Keane J., Mecum B., Moldovan-Grünfeld D., Ooms J., Wujciak-Jens K., Luraschi J., Dunkle Werner K., Wong J., Apache Arrow *arrow: Integration to 'Apache' 'Arrow'* <https://cran.r-project.org/web/packages/arrow/index.html> (2025).
54. Mikryukov, V. *vmikk/metagMisc: miscellaneous functions for metagenomic analysis*. <https://github.com/vmikk/metagMisc> (2025).
55. Chiu, C.-H. & Chao, A. Estimating and comparing microbial diversity in the presence of sequencing errors. *PeerJ* **4**, e1634 (2016).
56. Jost, L. Entropy and diversity. *Oikos* **113**, 363–375 (2006).

57. Pebesma, E., Bivand R., Racine E., Sumner M., Cook I., Keitt T., Lovelace R., Wickham H., Ooms J., Müller K., Pedersen TL., Baston D., Dunnington D. *sf: Simple Features for R* <https://cran.r-project.org/web/packages/sf/index.html> (2025).
58. Šavrič, B., Patterson, T. & Jenny, B. The Equal Earth map projection. *Int. J. Geogr. Inf. Sci.* **33**, 454–465 (2019).
59. GADM development team. Global Administrative Areas GADM version 2.8 <https://gadm.org/index.html> (2025).
60. Loidi, J., Navarro-Sánchez, G. & Vynokurov, D. A vector map of the world's terrestrial biotic units: subbiomes, biomes, ecozones and domains. *Veg. Classif. Surv.* **4**, 59–61 (2023).
61. Dinerstein, E. *et al.* An Ecoregion-Based Approach to Protecting Half the Terrestrial Realm. *Bioscience* **67**, 534–545 (2017).
62. FAO; International Institute for Applied Systems Analysis (IIASA). *Harmonized World Soil Database Version 2.0*. doi:10.4060/cc3823en (2023).
63. Meyer, H. & Pebesma, E. Predicting into unknown space? Estimating the area of applicability of spatial prediction models. *Methods Ecol. Evol.* **12**, 1620–1633 (2021).
64. Carter M.R., Gregorich E.G., *Soil Sampling and Methods of Analysis*. (CRC Press, Boca Raton, 2007). doi:10.1201/9781420005271.
65. Jurburg, S. D., Keil, P., Singh, B. K. & Chase, J. M. All together now: Limitations and recommendations for the simultaneous analysis of all eukaryotic soil sequences. *Mol. Ecol. Resour.* **21**, 1759–1771 (2021).
66. Chen, M. *et al.* Sampling Design and Sample Processing Affect Soil Biodiversity Assessments. *Mol. Ecol. Resour.* **26**, e70113 (2026).
67. Dopheide, A., Xie, D., Buckley, T. R., Drummond, A. J. & Newcomb, R. D. Impacts of DNA extraction and PCR on DNA metabarcoding estimates of soil biodiversity. *Methods Ecol. Evol.* **10**, 120–133 (2019).
68. Sun, T., Wang, Y., Lucas-Borja, M. E., Jing, X. & Feng, W. Divergent vertical distributions of microbial biomass with soil depth among groups and land uses. *J. Environ. Manage.* **292**, 112755 (2021).
69. Gueidan, C. *et al.* PacBio amplicon sequencing for metabarcoding of mixed DNA samples from lichen herbarium specimens. *MycoKeys* **53**, 73–91 (2019).
70. Franklin J., *Mapping Species Distributions: Spatial Inference and Prediction* (Cambridge University Press, Cambridge, 2012). <https://doi.org/10.1017/CBO9780511810602>.
71. Beck, H. E. *et al.* Present and future Köppen-Geiger climate classification maps at 1-km resolution. *Sci. Data* **5**, 180214 (2018).
72. Guerra, C. A. *et al.* Blind spots in global soil biodiversity and ecosystem function research. *Nat. Commun.* **11**, 3870 (2020).
73. Tedersoo, L. *et al.* Thirty novel fungal lineages: formal description based on environmental samples and DNA. *MycoKeys* **124**, 1–121 (2025).
74. Yatsiuk, I. *et al.* Arcyria and allied genera: taxonomic backbone and character evolution. *Fungal Syst. Evol.* **15**, 97–120 (2025).

75. Škaloud, P., Tučková, K., Čablová, R., Jadrná, I. & Černajová, I. High-frequency sampling unveils biotic and abiotic drivers of rapid phytoplankton morphological changes. *New Phytol.* **248**, 2528–2541 (2025).
76. Geisen, S. *et al.* Soil protists: a fertile frontier in soil biology research. *FEMS Microbiol. Rev.* **42**, 293–323 (2018).
77. Berlinches de Gea, A. *et al.* Protists as determinants of the One Health framework. *ISME J.* **19**, wraf179 (2025).
78. Derevnina, L. *et al.* Emerging oomycete threats to plants and animals. *Philos. Trans. R. Soc. B Biol. Sci.* **371**, 20150459 (2016).
79. Pölme, S. *et al.* FungalTraits: a user-friendly traits database of fungi and fungus-like stramenopiles. *Fungal Divers.* **105**, 1–16 (2020).
80. Nemaplex development team. *Nemaplex: The nematode information system.* (UCDavis.edu: Revision Date: 12/02/2025) <https://nemaplex.ucdavis.edu/index.htm> (2025).
81. Freudenthal, J., Schlegel, M., Bonkowski, M. & Dumack, K. A Novel Protistan Trait Database Reveals Functional Redundancy and Complementarity in Terrestrial Protists (Amoebozoa and Rhizaria). *Mol. Ecol. Resour.* **26**, e70064 (2026).
82. Kattge, J. *et al.* TRY plant trait database – enhanced coverage and open access. *Glob. Change Biol.* **26**, 119–188 (2020).
83. Hug, L. A. *et al.* A roadmap for equitable reuse of public microbiome data. *Nat. Microbiol.* **10**, 2384–2395 (2025).