

## Global mapping of groundwater-dependent ecosystems in drylands

### Technical documentation for:

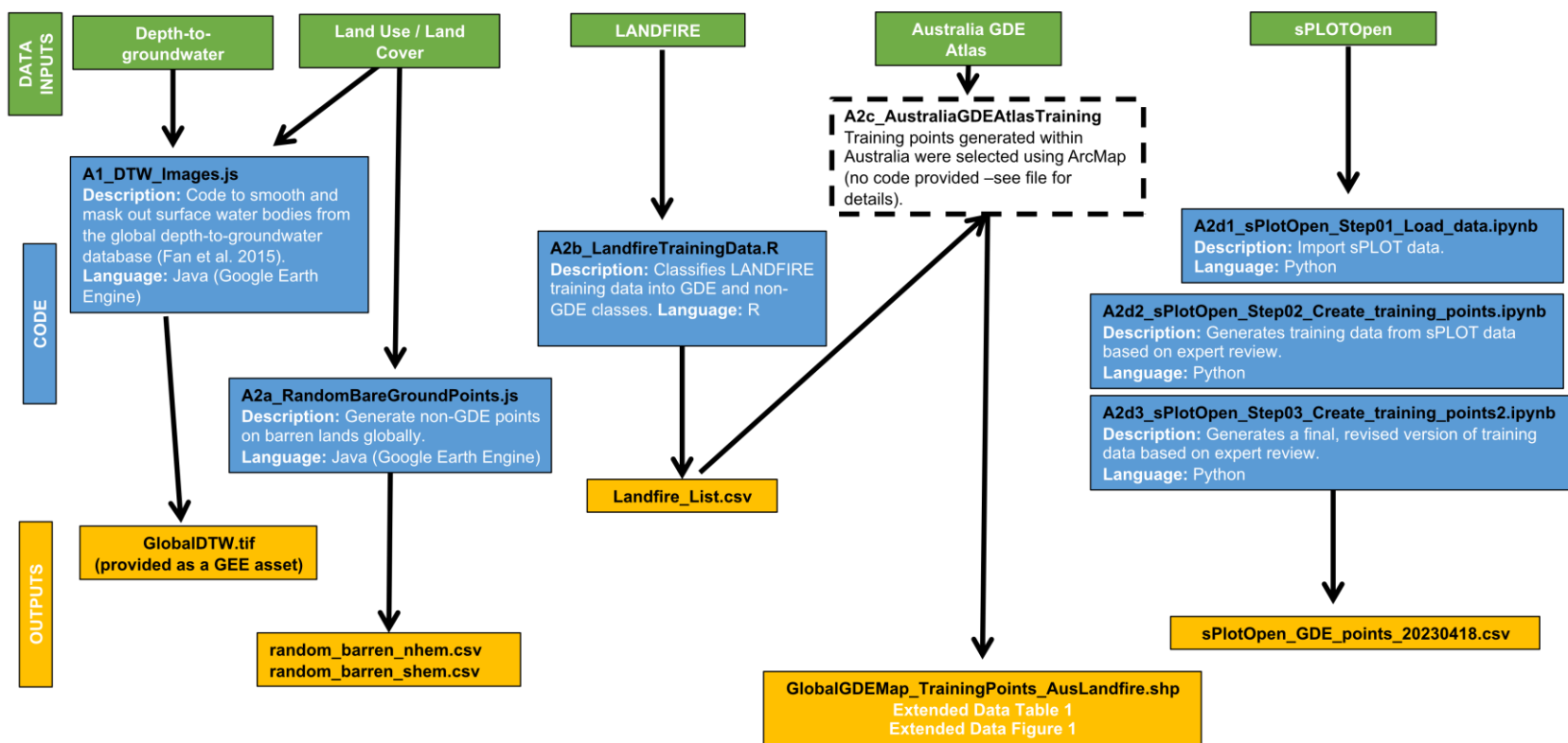
Rohde, M.M., C.M. Albano, X. Huggins, K.R. Klausmeyer, C. Morton, A. Sharman, E. Zaveri, L. Saito, Z. Freed, J.K. Howard, N. Job, H. Richter, K. Toderich, A. Rodella, T. Gleeson, J. Huntington, H.A. Chandanpurkar, A.J. Purdy, J.S. Famiglietti, M.B. Singer, D.A. Roberts, K. Caylor, J.C. Stella. Groundwater-dependent ecosystem map exposes global dryland protection needs. *Nature* (2024). DOI: [10.1038/s41586-024-07702-8](https://doi.org/10.1038/s41586-024-07702-8)

Prepared by Melissa M. Rohde and Xander Huggins  
[melissa@RohdeEnvironmental.com](mailto:melissa@RohdeEnvironmental.com)  
[xanderhuggins@uvic.ca](mailto:xanderhuggins@uvic.ca)

This technical documentation outlines the data and code used to map and analyze groundwater-dependent ecosystems across global drylands, as well as produce final data outputs. Each of the four sections contain a flow chart illustrating what data inputs (green boxes) were used to execute code (blue boxes), and the resulting outputs (gold boxes). Individual scripts are described in tables below each flow chart.

Section A GDE Model preparation	Section B GDE Random Forest model	Section C Post hoc analysis & result reporting	Section D Data access
<ul style="list-style-type: none"><li>Preprocesses input data for random forest model.</li><li>Generates training point dataset for the random forest model.</li></ul> <p>Output from these scripts feeds into section B.</p>	<p>Random forest (RF) model development including:</p> <ul style="list-style-type: none"><li>Assigning predictor variables to each training point for model ingestion.</li><li>Hyperparameter tuning.</li><li>Running RF model.</li><li>Model validation: Cross Validation tests.</li><li>Model validation: distribution plots of training data.</li></ul> <p>Generates core output of this study: <b>Global 30m GDE classification map and GDE probability map.</b></p>	<ul style="list-style-type: none"><li>Prepares groundwater storage and protected area data for comparison to GDE map.</li><li>Assesses GDE distributions against these post hoc datasets.</li><li>Derives all quantitative results.</li><li>Prepares aggregated data products for deposition (see Section D).</li></ul>	<ul style="list-style-type: none"><li>Description of all deposited data products:<ul style="list-style-type: none"><li>- 1 arcsec (~30 m),</li><li>- 30 arcsec (~1 km),</li><li>- 5 arcmin (~10 km),</li><li>- 30 arcmin (~50 km).</li></ul></li><li>Data access URLs.</li></ul>
<b>PAGE 2</b> <a href="#">[go to section]</a>	<b>PAGE 5</b> <a href="#">[go to section]</a>	<b>PAGE 20</b> <a href="#">[go to section]</a>	<b>PAGE 25</b> <a href="#">[go to section]</a>

## Section A GDE Model Preparation



**DATA INPUTS:** Data inputs are provided in the Zenodo repo. The original data inputs used are based on publicly available datasets that need to be directly downloaded from the weblinks included within the code.

**CODE:** All code was performed using the Google Earth Engine code editor (Java), R version 4.3.1, or Python 3.9.15

**OUTPUT:** Output data files are provided either in the Zenodo repository and/or as a google earth engine (GEE) asset.

## A1 – Prepare masking data

<b>Script name</b>	A1_DTG_Images.js
<b>Description</b>	Smooths and masks out surface water bodies from the global depth-to-groundwater database (Fan et al. 2017)
<b>Data inputs</b>	Fan et al. 2017 (subset by continent; Google Earth Engine asset) ESRI 10 m Land Use Land Cover data (Google Earth Engine asset)
<b>Intermediary data outputs</b>	GlobalDTG_1.tif, GlobalDTG_2.tif (provided as a Google Earth Engine assets)
<b>Authors</b>	Christine Albano and Melissa M. Rohde

## A2 – Generate Training Data

<b>Script name</b>	A2a_RandomBareGroundPoints.js
<b>Description</b>	Generates 10,000 non-GDE points on barren lands globally (5000 points in each hemisphere).
<b>Data inputs</b>	ESRI 10 m Land Use Land Cover data (Google Earth Engine asset)
<b>Intermediary data outputs</b>	random_barren_nhem.csv random_barren_shem.csv
<b>Authors</b>	Christine Albano and Melissa M. Rohde

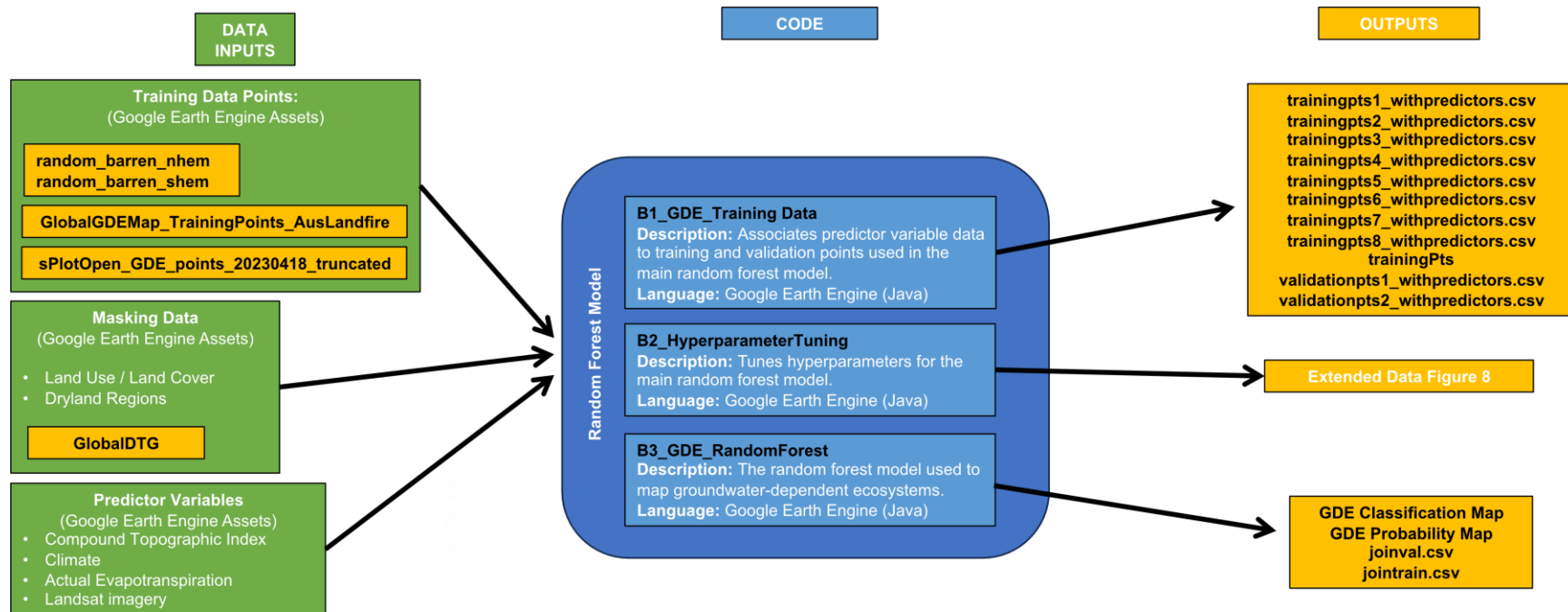
<b>Script name</b>	A2b_LandfireTrainingData.R
<b>Description</b>	Classifies LANDFIRE training data into GDE and non-GDE classes.
<b>Data inputs</b>	vegcamp_GDE_review_20161107.csv: This file contains GDE certainty from California based on Klausmeyer et al. 2018.  Landfire_TrainingData.csv: This file contains GDE confirmation based on published data and expert review (see Supplementary Table 4).  tbl_DomSp_Unique_NW_SW.csv: This file contains a list of all the unique dominant species names in the LANDFIRE dataset.
<b>R packages used</b>	<i>dplyr</i> (Wickham et al., 2023).
<b>Intermediary data outputs</b>	Landfire_List.csv: This file was used to extract training data from the original LANDFIRE reference data in ArcGIS. GDE classification for individual states is as follows: 0=N/A, 1=GDE, 2 = non-GDE. GDE certainty ranges from 1-3, with 1 being the highest, and 3 being the lowest.

<b>Script name</b>	A2c_AustraliaGDEAtlasTraining.rtf
<b>Description</b>	Training points generated within Australia were selected using ArcMap (no code provided). These are the steps for processing the Australian GDE training points.
<b>Data inputs</b>	GDE_Atlas_Aquatic_GDEs.gdb ( $n = 1,107,524$ features) GDE_Atlas_Terrestrial_GDEs.gdb ( $n = 7,747,955$ features)  Both geodatabases are not provided in this data repository. Contact: <a href="mailto:water@bom.gov.au">water@bom.gov.au</a> at the Australian Bureau of Meteorology for bulk download information.
<b>Intermediary data outputs</b>	GlobalGDEMap_TrainingPoints_AusLandfire.shp
<b>Authors</b>	Melissa M. Rohde

<b>Script name</b>	A2d1_sPlotOpen_Step01_Load_data.ipynb A2d2_sPlotOpen_Step02_Create_training_points.ipynb A2d3_sPlotOpen_Step03_Create_training_points2.ipynb
<b>Description</b>	The first python code file includes the steps to import the sPLOT Open dataset and summarize all of the species found in each continent. This table was then shared with expert reviewers to classify each species/continent combination as GDE or not GDE. The second code file takes the results of the expert review and links it back to the sPLOT Open dataset to create the first version of the training points for the machine learning model. The third code file imports a revised version of the expert review file and creates the final set of training points.
<b>Data inputs</b>	sPlotOpen_GDE_review_20230407.csv  sPlotOpen_GDE_review_20230407_forToderichReview_Kristina Last_KK.xlsx Both of the above are developed from the sPlotOpen dataset (Sabatini et al., 2021).
<b>Python modules used</b>	<i>Pandas</i> (Pandas development team, 2020; McKinney, 2010) <i>NumPy</i> (Harris et al., 2020). <i>FuzzyWuzzy</i> (Cohen, 2020).
<b>Intermediary data outputs</b>	sPlotOpen_GDE_points_20230418_truncated.csv sPlotOpen_FinalList_20230418.csv sPLOTOpen_GDE_points_20230418.csv
<b>Authors</b>	Kirk Klausmeyer

## Section B GDE Random Forest Model

### SUBSECTION: MODEL DEVELOPMENT



**DATA INPUTS:** Data inputs are provided as Google Earth Engine assets within the code.

**CODE:** All code was performed using the Google Earth Engine code editor (Java).

**OUTPUT:** Output data files (black text) are provided either in the Zenodo repository and/or as a google earth engine asset, or those in white text are provided in the published paper.

## B1 – Export Training Data for main model

<b>Script name</b>	B1_GDE_TrainingData.js
<b>Description</b>	Associates predictor variable data to training and validation points used in the B3_GDE_RandomForest mapping model.
<b>Data inputs</b>	<p>All data inputs are uploaded and accessible as Google Earth Engine assets. Links are provided in the code.</p> <p>Masking Data:</p> <ul style="list-style-type: none"> <li>• GlobalDTG</li> <li>• ESRI 10 m Land Use Land Cover data</li> <li>• Dryland Regions (Beck_KG_V1_present_0p0083)</li> </ul> <p>Predictor Variables:</p> <ul style="list-style-type: none"> <li>• Compound Topographic Index (TopConvIndex_Global_KGclims)</li> <li>• Climate - precipitation and potential evapotranspiration (IDAHO_EPSCOR/TERRACLIMATE)</li> <li>• Actual Evapotranspiration (CAS/IGSNRR/PML/V2_v017)</li> <li>• Landsat imagery (LANDSAT/LC08/C02/T1_L2)</li> </ul> <p>Training Data Points:</p> <ul style="list-style-type: none"> <li>• GlobalGDEMap_TrainingPoints_AusLandfire</li> <li>• random_barren_nhem</li> <li>• random_barren_shem</li> <li>• sPlotOpen_GDE_points_20230418_truncated</li> </ul>
<b>Intermediary data outputs</b>	trainingpts1_withpredictors.csv trainingpts2_withpredictors.csv trainingpts3_withpredictors.csv trainingpts4_withpredictors.csv trainingpts5_withpredictors.csv trainingpts6_withpredictors.csv trainingpts7_withpredictors.csv trainingpts8_withpredictors.csv trainingPts (Google Earth Engine asset) validationpts1_withpredictors.csv validationpts2_withpredictors.csv
<b>Authors</b>	Melissa M. Rohde

## B2 - Hyperparameter Tuning

<b>Script name</b>	B2_HyperparameterTuning.js
<b>Description</b>	Retrieves optimal parameters in the random forest model to achieve highest accuracy for GDE classification.
<b>Data inputs</b>	<p>All data inputs are uploaded and accessible as Google Earth Engine assets. Links are provided in the code.</p> <p>Masking Data:</p> <ul style="list-style-type: none"> <li>• GlobalDTG</li> <li>• ESRI 10 m Land Use Land Cover data</li> <li>• Dryland Regions (Beck_KG_V1_present_0p0083)</li> </ul> <p>Predictor Variables:</p> <ul style="list-style-type: none"> <li>• Compound Topographic Index (TopConvIndex_Global_KGclims)</li> <li>• Climate - precipitation and potential evapotranspiration (IDAHO_EPSCOR/TERRACLIMATE)</li> <li>• Actual Evapotranspiration (CAS/IGSNRR/PML/V2_v017)</li> <li>• Landsat imagery (LANDSAT/LC08/C02/T1_L2)</li> </ul> <p>Training Data:</p> <ul style="list-style-type: none"> <li>• trainingpts1_withpredictors</li> <li>• trainingpts2_withpredictors</li> <li>• trainingpts3_withpredictors</li> <li>• trainingpts4_withpredictors</li> <li>• trainingpts5_withpredictors</li> <li>• trainingpts6_withpredictors</li> <li>• trainingpts7_withpredictors</li> <li>• trainingpts8_withpredictors</li> <li>• validationpts1_withpredictors</li> <li>• validationpts2_withpredictors</li> </ul>
<b>Intermediary data outputs</b>	Accuracy tables and charts for <i>numberOfTrees</i> , <i>variablesPerSplit</i> , <i>minLeafPopulation</i> , <i>bagFraction</i> , and <i>maxNodes</i> parameters. Output is compiled in Extended Data Figure 8.
<b>Authors</b>	Melissa M. Rohde

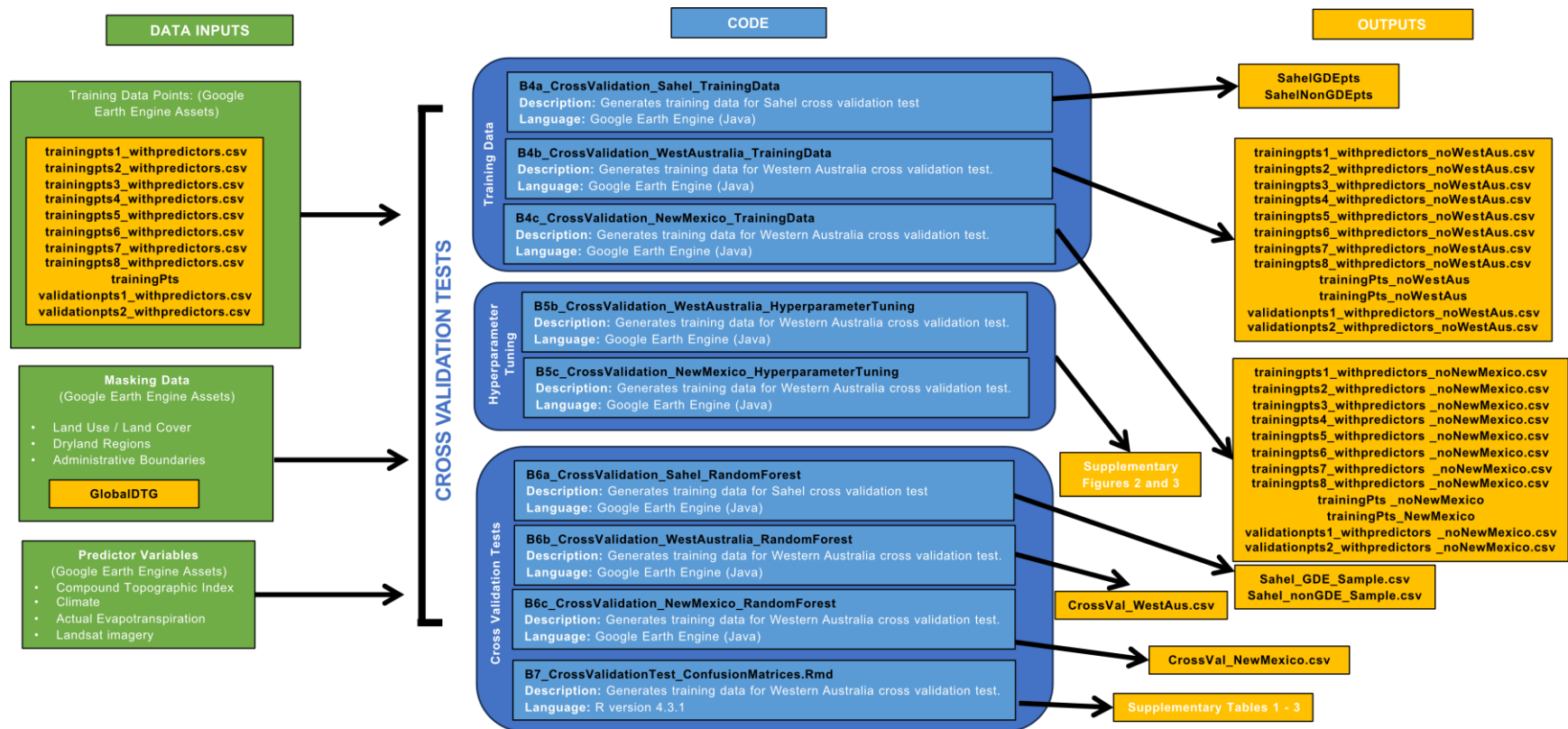
### B3 - Classify GDEs using the Random Forest model

<b>Script name</b>	B3_GDE_RandomForest.js
<b>Description</b>	The random forest model used to map groundwater-dependent ecosystems globally. Tuned hyperparameters from B2 (above) are included in this model run.
<b>Data inputs</b>	<p>All data inputs are uploaded and accessible as Google Earth Engine assets. Links are provided in the code.</p> <p>Masking Data:</p> <ul style="list-style-type: none"> <li>• GlobalDTG</li> <li>• ESRI 10 m Land Use Land Cover data</li> <li>• Dryland Regions (Beck_KG_V1_present_0p0083)</li> </ul> <p>Predictor Variables:</p> <ul style="list-style-type: none"> <li>• Compound Topographic Index (TopConvIndex_Global_KGclims)</li> <li>• Climate - precipitation and potential evapotranspiration (IDAHO_EPSCOR/TERRACLIMATE)</li> <li>• Actual Evapotranspiration (CAS/IGSNRR/PML/V2_v017)</li> <li>• Landsat imagery (LANDSAT/LC08/C02/T1_L2)</li> </ul> <p>Training Data:</p> <ul style="list-style-type: none"> <li>• trainingpts1_withpredictors</li> <li>• trainingpts2_withpredictors</li> <li>• trainingpts3_withpredictors</li> <li>• trainingpts4_withpredictors</li> <li>• trainingpts5_withpredictors</li> <li>• trainingpts6_withpredictors</li> <li>• trainingpts7_withpredictors</li> <li>• trainingpts8_withpredictors</li> <li>• validationpts1_withpredictors</li> <li>• validationpts2_withpredictors</li> </ul>
<b>Intermediary data outputs</b>	<p>GDE classification map (Google Earth Engine asset)</p> <p>GDE probability map (Google Earth Engine asset)</p> <p>joinval.csv</p> <p>jointrain.csv</p>
<b>Authors</b>	Christine Albano and Melissa M. Rohde



## Section B: GDE Random Forest Model

### SUBSECTION: MODEL VALIDATION - CROSS VALIDATION



**DATA INPUTS:** Data inputs are provided as Google Earth Engine assets within the code.

**CODE:** All code was performed using the Google Earth Engine code editor (Java) or R version 4.3.1.

**OUTPUT:** Output data files (black text) are provided either in the Zenodo repository and/or as a Google Earth Engine asset, or those in white text are provided in the published paper.

#### B4 - Generate training data for regional cross validation tests

<b>Script name</b>	B4a_CrossValidation_Sahel_TrainingData.js
<b>Description</b>	Generates training data points using GDE location data from the World Bank that were extracted from peer-review literature sources, and non-GDE points from the ESRI land use and land cover bare ground layer.
<b>Data inputs</b>	GDE lines and point data ( <a href="#">Source: Rodella et al. 2023</a> ) ESRI 10 m Land Use Land Cover data (Google Earth Engine asset)
<b>Intermediary data outputs</b>	SahelGDEpts (Google Earth Engine asset) SahelNonGDEpts (Google Earth Engine asset)
<b>Authors</b>	Melissa M. Rohde

<b>Script name</b>	B4b_CrossValidation_WestAustralia_TrainingData.js
<b>Description</b>	This Google Earth Engine code generates training data points for the Western Australia cross validation test.
<b>Data inputs</b>	<p>All data inputs are uploaded and accessible as Google Earth Engine assets. Links are provided in the code.</p> <p>Masking Data:</p> <ul style="list-style-type: none"> <li>• GlobalDTG</li> <li>• ESRI 10 m Land Use Land Cover data</li> <li>• Dryland Regions (Beck_KG_V1_present_0p0083)</li> </ul> <p>Predictor Variables:</p> <ul style="list-style-type: none"> <li>• Compound Topographic Index (TopConvIndex_Global_KGclims)</li> <li>• Climate - precipitation and potential evapotranspiration (IDAHO_EPSCOR/TERRACLIMATE)</li> <li>• Actual Evapotranspiration (CAS/IGSNRR/PML/V2_v017)</li> <li>• Landsat imagery (LANDSAT/LC08/C02/T1_L2)</li> </ul> <p>Training Data Points:</p> <ul style="list-style-type: none"> <li>• Administrative boundaries (FAO/GAUL/2015/level1)</li> <li>• GlobalGDEMap_TrainingPoints_AusLandfire</li> <li>• random_barren_nhem</li> <li>• random_barren_shem</li> <li>• sPlotOpen_GDE_points_20230418_truncated</li> </ul>
<b>Intermediary data outputs</b>	trainingpts1_withpredictors_noWestAus.csv trainingpts2_withpredictors_noWestAus.csv trainingpts3_withpredictors_noWestAus.csv trainingpts4_withpredictors_noWestAus.csv trainingpts5_withpredictors_noWestAus.csv trainingpts6_withpredictors_noWestAus.csv trainingpts7_withpredictors_noWestAus.csv trainingpts8_withpredictors_noWestAus.csv trainingPts_noWestAus (Google Earth Engine asset) trainingPts_WestAus (Google Earth Engine asset)

	validationpts1_withpredictors_noWestAus.csv validationpts2_withpredictors_noWestAus.csv
<b>Authors</b>	Melissa M. Rohde

<b>Script name</b>	B4c_CrossValidation_NewMexico_TrainingData.js
<b>Description</b>	Generates training data points for the New Mexico cross validation test.
<b>Data inputs</b>	<p>All data inputs are uploaded and accessible as Google Earth Engine assets. Links are provided in the code.</p> <p>Masking Data:</p> <ul style="list-style-type: none"> <li>• GlobalDTG</li> <li>• ESRI 10 m Land Use Land Cover data</li> <li>• Dryland Regions (Beck_KG_V1_present_0p0083)</li> </ul> <p>Predictor Variables:</p> <ul style="list-style-type: none"> <li>• Compound Topographic Index (TopConvIndex_Global_KGclims)</li> <li>• Climate - precipitation and potential evapotranspiration (IDAHO_EPSCOR/TERRACLIMATE)</li> <li>• Actual Evapotranspiration (CAS/IGSNRR/PML/V2_v017)</li> <li>• Landsat imagery (LANDSAT/LC08/C02/T1_L2)</li> </ul> <p>Training Data Points:</p> <ul style="list-style-type: none"> <li>• Administrative boundaries (FAO/GAUL/2015/level1)</li> <li>• GlobalGDEMap_TrainingPoints_AusLandfire</li> <li>• random_barren_nhem</li> <li>• random_barren_shem</li> <li>• sPlotOpen_GDE_points_20230418_truncated</li> </ul>
<b>Intermediary data outputs</b>	trainingpts1_withpredictors_noNewMexico.csv trainingpts2_withpredictors_noNewMexico.csv trainingpts3_withpredictors_noNewMexico.csv trainingpts4_withpredictors_noNewMexico.csv trainingpts5_withpredictors_noNewMexico.csv trainingpts6_withpredictors_noNewMexico.csv trainingpts7_withpredictors_noNewMexico.csv trainingpts8_withpredictors_noNewMexico.csv trainingPts_noNewMexico (Google Earth Engine asset) trainingPts_NewMexico (Google Earth Engine asset) validationpts1_withpredictors_noNewMexico.csv validationpts2_withpredictors_noNewMexico.csv
<b>Authors</b>	Melissa M. Rohde

**B5 - Model Validation: Cross Validation Tests. Hyperparameter tuning for regional cross validation tests**

<b>Script name</b>	B2_HyperparameterTuning.js – name doesn't match section?
<b>Description</b>	Retrieves optimal parameters in random forest model to achieve highest accuracy for GDE classification in the Sahel cross validation test.
<b>Data inputs</b>	<p>All data inputs are uploaded and accessible as Google Earth Engine assets. Links are provided in the code.</p> <p>Masking Data:</p> <ul style="list-style-type: none"> <li>GlobalDTG</li> <li>ESRI 10 m Land Use Land Cover data</li> <li>Dryland Regions (Beck_KG_V1_present_0p0083)</li> </ul> <p>Predictor Variables:</p> <ul style="list-style-type: none"> <li>Compound Topographic Index (TopConvIndex_Global_KGclims)</li> <li>Climate - precipitation and potential evapotranspiration (IDAHO_EPSCOR/TERRACLIMATE)</li> <li>Actual Evapotranspiration (CAS/IGSNRR/PML/V2_v017)</li> <li>Landsat imagery (LANDSAT/LC08/C02/T1_L2)</li> </ul> <p>Training Data:</p> <ul style="list-style-type: none"> <li>trainingpts1_withpredictors</li> <li>trainingpts2_withpredictors</li> <li>trainingpts3_withpredictors</li> <li>trainingpts4_withpredictors</li> <li>trainingpts5_withpredictors</li> <li>trainingpts6_withpredictors</li> <li>trainingpts7_withpredictors</li> <li>trainingpts8_withpredictors</li> <li>validationpts1_withpredictors</li> <li>validationpts2_withpredictors</li> </ul>
<b>Intermediary data outputs</b>	Accuracy tables and charts for <i>numberOfTrees</i> , <i>variablesPerSplit</i> , <i>minLeafPopulation</i> , <i>bagFraction</i> , and <i>maxNodes</i> parameters.
<b>Authors</b>	Melissa M. Rohde

<b>Script name</b>	B5b_CrossValidation_WestAustralia_HyperparameterTuning.js
<b>Description</b>	Tunes parameters in the random forest model for the West Australia cross validation test to retrieve optimal parameters to achieve highest accuracy.
<b>Data inputs</b>	<p>All data inputs are uploaded and accessible as Google Earth Engine assets. Links are provided in the code.</p> <p>Masking Data:</p> <ul style="list-style-type: none"> <li>GlobalDTG</li> <li>ESRI 10 m Land Use Land Cover data</li> <li>Dryland Regions (Beck_KG_V1_present_0p0083)</li> </ul> <p>Predictor Variables:</p>

	<ul style="list-style-type: none"> <li>• Compound Topographic Index (TopConvIndex_Global_KGclims)</li> <li>• Climate - precipitation and potential evapotranspiration (IDAHO_EPSCOR/TERRACLIMATE)</li> <li>• Actual Evapotranspiration (CAS/IGSNRR/PML/V2_v017)</li> <li>• Landsat imagery (LANDSAT/LC08/C02/T1_L2)</li> </ul> <p>Training Data:</p> <ul style="list-style-type: none"> <li>• trainingpts1_withpredictors_noWestAus</li> <li>• trainingpts2_withpredictors_noWestAus</li> <li>• trainingpts3_withpredictors_noWestAus</li> <li>• trainingpts4_withpredictors_noWestAus</li> <li>• trainingpts5_withpredictors_noWestAus</li> <li>• trainingpts6_withpredictors_noWestAus</li> <li>• trainingpts7_withpredictors_noWestAus</li> <li>• trainingpts8_withpredictors_noWestAus</li> <li>• trainingPts_noWestAus</li> <li>• trainingPts_WestAus</li> <li>• validationpts1_withpredictors_noWestAus</li> <li>• validationpts2_withpredictors_noWestAus</li> </ul>
<b>Intermediary data outputs</b>	Accuracy tables and charts for <i>numberOfTrees</i> , <i>variablesPerSplit</i> , <i>minLeafPopulation</i> , <i>bagFraction</i> , and <i>maxNodes</i> parameters.
<b>Authors</b>	Melissa M. Rohde

<b>Script name</b>	B5c_CrossValidation_NewMexico_HyperparameterTuning.js
<b>Description</b>	Tunes parameters in the random forest model for the New Mexico cross validation test to retrieve optimal parameters to achieve highest accuracy.
<b>Data inputs</b>	<p>All data inputs are uploaded and accessible as Google Earth Engine assets. Links are provided in the code.</p> <p>Masking Data:</p> <ul style="list-style-type: none"> <li>• GlobalDTG</li> <li>• ESRI 10 m Land Use Land Cover data</li> <li>• Dryland Regions (Beck_KG_V1_present_0p0083)</li> </ul> <p>Predictor Variables:</p> <ul style="list-style-type: none"> <li>• Compound Topographic Index (TopConvIndex_Global_KGclims)</li> <li>• Climate - precipitation and potential evapotranspiration (IDAHO_EPSCOR/TERRACLIMATE)</li> <li>• Actual Evapotranspiration (CAS/IGSNRR/PML/V2_v017)</li> <li>• Landsat imagery (LANDSAT/LC08/C02/T1_L2)</li> </ul> <p>Training Data:</p> <ul style="list-style-type: none"> <li>• trainingpts1_withpredictors_noNewMexico</li> <li>• trainingpts2_withpredictors_noNewMexico</li> <li>• trainingpts3_withpredictors_noNewMexico</li> <li>• trainingpts4_withpredictors_noNewMexico</li> <li>• trainingpts5_withpredictors_noNewMexico</li> <li>• trainingpts6_withpredictors_noNewMexico</li> <li>• trainingpts7_withpredictors_noNewMexico</li> <li>• trainingpts8_withpredictors_noNewMexico</li> <li>• trainingPts_noNewMexico</li> <li>• trainingPts_NewMexico</li> <li>• validationpts1_withpredictors_noNewMexico</li> </ul>

	<ul style="list-style-type: none"> <li>validationpts2_withpredictors_noNewMexico</li> </ul>
<b>Intermediary data outputs</b>	Accuracy tables and charts for <i>numberOfTrees</i> , <i>variablesPerSplit</i> , <i>minLeafPopulation</i> , <i>bagFraction</i> , and <i>maxNodes</i> parameters.
<b>Authors</b>	Melissa M. Rohde

**B6 - Model Validation: Cross Validation Tests. Perform regional cross validation tests with tuned hyperparameters**

<b>Script name</b>	B6a_CrossValidation_Sahel_RandomForest.js
<b>Description</b>	Samples the main global GDE classification output using the Sahel training points.
<b>Data inputs</b>	<p>All data inputs are uploaded and accessible as Google Earth Engine assets. Links are provided in the code.</p> <p>Masking Data:</p> <ul style="list-style-type: none"> <li>GlobalDTG</li> <li>ESRI 10 m Land Use Land Cover data</li> <li>Dryland Regions (Beck_KG_V1_present_0p0083)</li> </ul> <p>Predictor Variables:</p> <ul style="list-style-type: none"> <li>Compound Topographic Index (TopConvIndex_Global_KGclims)</li> <li>Climate - precipitation and potential evapotranspiration (IDAHO_EPSCOR/TERRACLIMATE)</li> <li>Actual Evapotranspiration (CAS/IGSNRR/PML/V2_v017)</li> <li>Landsat imagery (LANDSAT/LC08/C02/T1_L2)</li> </ul> <p>Training Data:</p> <ul style="list-style-type: none"> <li>trainingpts1_withpredictors</li> <li>trainingpts2_withpredictors</li> <li>trainingpts3_withpredictors</li> <li>trainingpts4_withpredictors</li> <li>trainingpts5_withpredictors</li> <li>trainingpts6_withpredictors</li> <li>trainingpts7_withpredictors</li> <li>trainingpts8_withpredictors</li> <li>validationpts1_withpredictors</li> <li>validationpts2_withpredictors</li> <li>trainingPts</li> <li>SahelGDEpts</li> <li>SahelNonGDEpts</li> </ul>
<b>Intermediary data outputs</b>	Sahel_GDE_Sample.csv Sahel_nonGDE_Sample.csv
<b>Authors</b>	Melissa M. Rohde

<b>Script name</b>	B6b_CrossValidation_WestAustralia_RandomForest.js
<b>Description</b>	Cross validation test for Western Australia.
<b>Data inputs</b>	<p>All data inputs are uploaded and accessible as Google Earth Engine assets. Links are provided in the code.</p> <p>Masking Data:</p> <ul style="list-style-type: none"> <li>• GlobalDTG</li> <li>• ESRI 10 m Land Use Land Cover data</li> <li>• Dryland Regions (Beck_KG_V1_present_0p0083)</li> </ul> <p>Predictor Variables:</p> <ul style="list-style-type: none"> <li>• Compound Topographic Index (TopConvIndex_Global_KGclims)</li> <li>• Climate - precipitation and potential evapotranspiration (IDAHO_EPSCOR/TERRACLIMATE)</li> <li>• Actual Evapotranspiration (CAS/IGSNRR/PML/V2_v017)</li> <li>• Landsat imagery (LANDSAT/LC08/C02/T1_L2)</li> </ul> <p>Training Data</p> <ul style="list-style-type: none"> <li>• trainingpts1_withpredictors_noWestAus</li> <li>• trainingpts2_withpredictors_noWestAus</li> <li>• trainingpts3_withpredictors_noWestAus</li> <li>• trainingpts4_withpredictors_noWestAus</li> <li>• trainingpts5_withpredictors_noWestAus</li> <li>• trainingpts6_withpredictors_noWestAus</li> <li>• trainingpts7_withpredictors_noWestAus</li> <li>• trainingpts8_withpredictors_noWestAus</li> <li>• trainingPts_noWestAus</li> <li>• trainingPts_WestAus</li> <li>• validationpts1_withpredictors_noWestAus</li> <li>• validationpts2_withpredictors_noWestAus</li> </ul>
<b>Intermediary data outputs</b>	CrossVal_WestAus.csv
<b>Authors</b>	Melissa M. Rohde

<b>Script name</b>	B6c_CrossValidation_NewMexico_HyperparameterTuning.js
<b>Description</b>	Cross validation test for New Mexico.
<b>Data inputs</b>	<p>All data inputs are uploaded and accessible as Google Earth Engine assets. Links are provided in the code.</p> <p>Masking Data:</p> <ul style="list-style-type: none"> <li>• GlobalDTG</li> <li>• ESRI 10 m Land Use Land Cover data</li> <li>• Dryland Regions (Beck_KG_V1_present_0p0083)</li> </ul> <p>Predictor Variables:</p> <ul style="list-style-type: none"> <li>• Compound Topographic Index (TopConvIndex_Global_KGclims)</li> <li>• Climate - precipitation and potential evapotranspiration (IDAHO_EPSCOR/TERRACLIMATE)</li> <li>• Actual Evapotranspiration (CAS/IGSNRR/PML/V2_v017)</li> <li>• Landsat imagery (LANDSAT/LC08/C02/T1_L2)</li> </ul> <p>Training Data:</p> <ul style="list-style-type: none"> <li>• trainingpts1_withpredictors_noNewMexico</li> <li>• trainingpts2_withpredictors_noNewMexico</li> </ul>

	<ul style="list-style-type: none"> <li>• trainingpts3_withpredictors_noNewMexico</li> <li>• trainingpts4_withpredictors_noNewMexico</li> <li>• trainingpts5_withpredictors_noNewMexico</li> <li>• trainingpts6_withpredictors_noNewMexico</li> <li>• trainingpts7_withpredictors_noNewMexico</li> <li>• trainingpts8_withpredictors_noNewMexico</li> <li>• trainingPts_noNewMexico</li> <li>• trainingPts_noNewMexico</li> <li>• validationpts1_withpredictors_noNewMexico</li> <li>• validationpts2_withpredictors_noNewMexico</li> </ul>
<b>Intermediary data outputs</b>	CrossVal_NewMexico.csv
<b>Authors</b>	Melissa M. Rohde

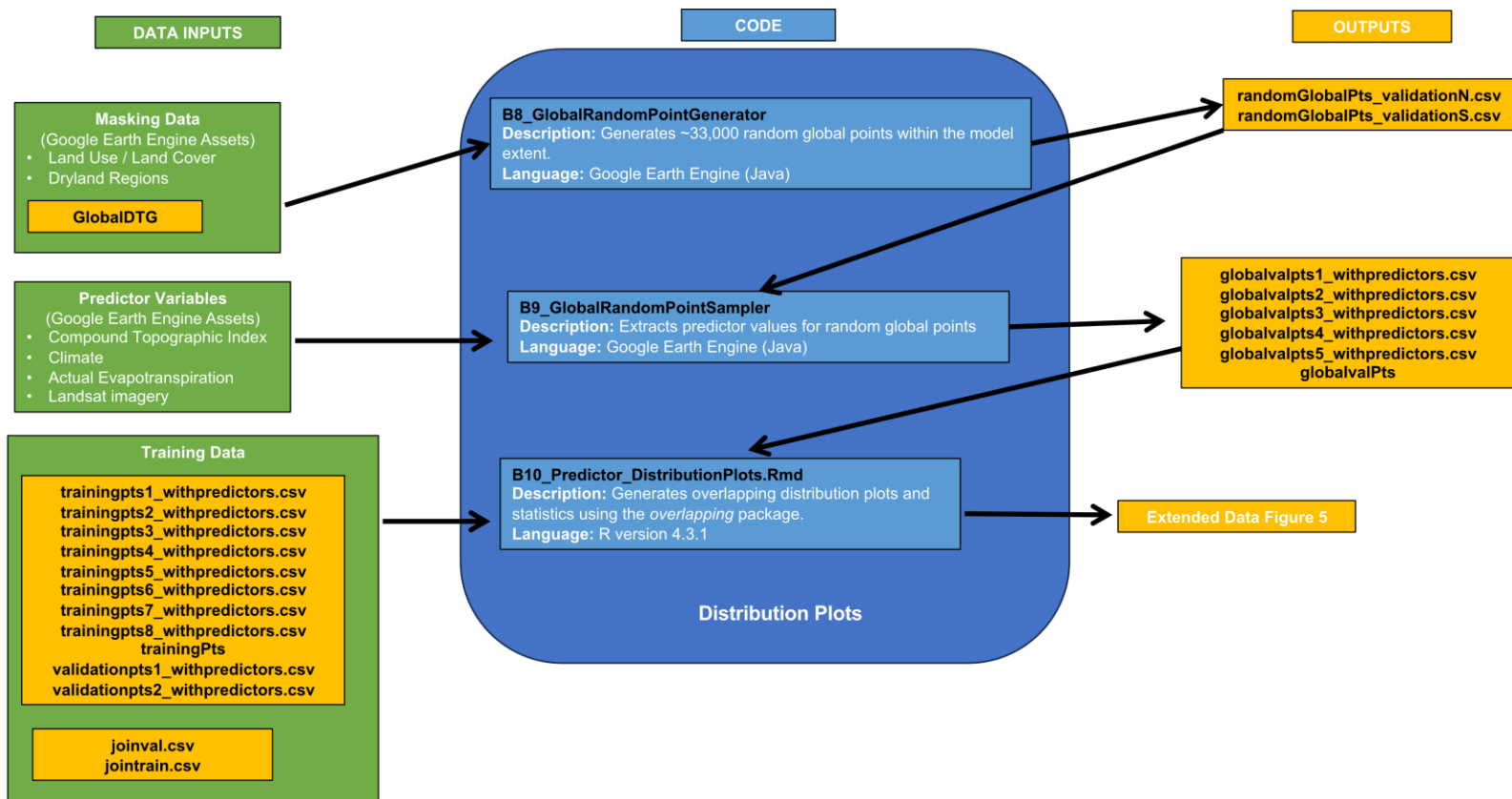
**B7 - Model Validation: Cross Validation Tests. Generate confusion matrices for cross validation tests**

<b>Script name</b>	B7_CrossValidationTests_ConfusionMatrices.Rmd
<b>Description</b>	Generate confusion matrices for the Sahel, Western Australia, and New Mexico cross validation tests.
<b>Data inputs</b>	Sahel_GDE_Sample.csv Sahel_nonGDE_Sample.csv CrossVal_WestAus.csv CrossVal_NewMexico.csv
<b>R packages used</b>	<i>overlapping</i> (Pastore et al., 2022). <i>gridExtra</i> (Auguie & Antonov, 2017). <i>tidyverse</i> (Wickham et al., 2019)
<b>Intermediary data outputs</b>	Confusion matrix for each cross-validation test (Supplementary Tables 1-3).



## Section B: GDE Random Forest Model

### SUBSECTION: MODEL VALIDATION – DISTRIBUTION PLOTTING



**DATA INPUTS:** Data inputs are provided as Google Earth Engine assets within the code or in the Zenodo repository.

**CODE:** All code was performed using the Google Earth Engine code editor (Java) or R version 4.3.1.

**OUTPUT:** Output data files (black text) are provided either in the Zenodo repository and/or as a google earth engine asset, or those in white text are provided in the published paper.

**B8 - Generate Random global points within model extent to compare distribution of variables within model extent and training data used in model.**

<b>Script name</b>	B8_GlobalRandomPointGenerator.js
<b>Description</b>	This Google Earth Engine code generates 125,000 random points, retaining only those that fall within the model extent (n~33,000 global points).
<b>Data inputs</b>	All data inputs are uploaded and accessible as Google Earth Engine assets. Links are provided in the code.  Masking Data: <ul style="list-style-type: none"> <li>• GlobalDTG</li> <li>• ESRI 10 m Land Use Land Cover data</li> <li>• Dryland Regions (Beck_KG_V1_present_0p0083)</li> </ul>
<b>Intermediary data outputs</b>	randomGlobalPts_validationN.csv randomGlobalPts_validationS.csv
<b>Authors</b>	Melissa M. Rohde

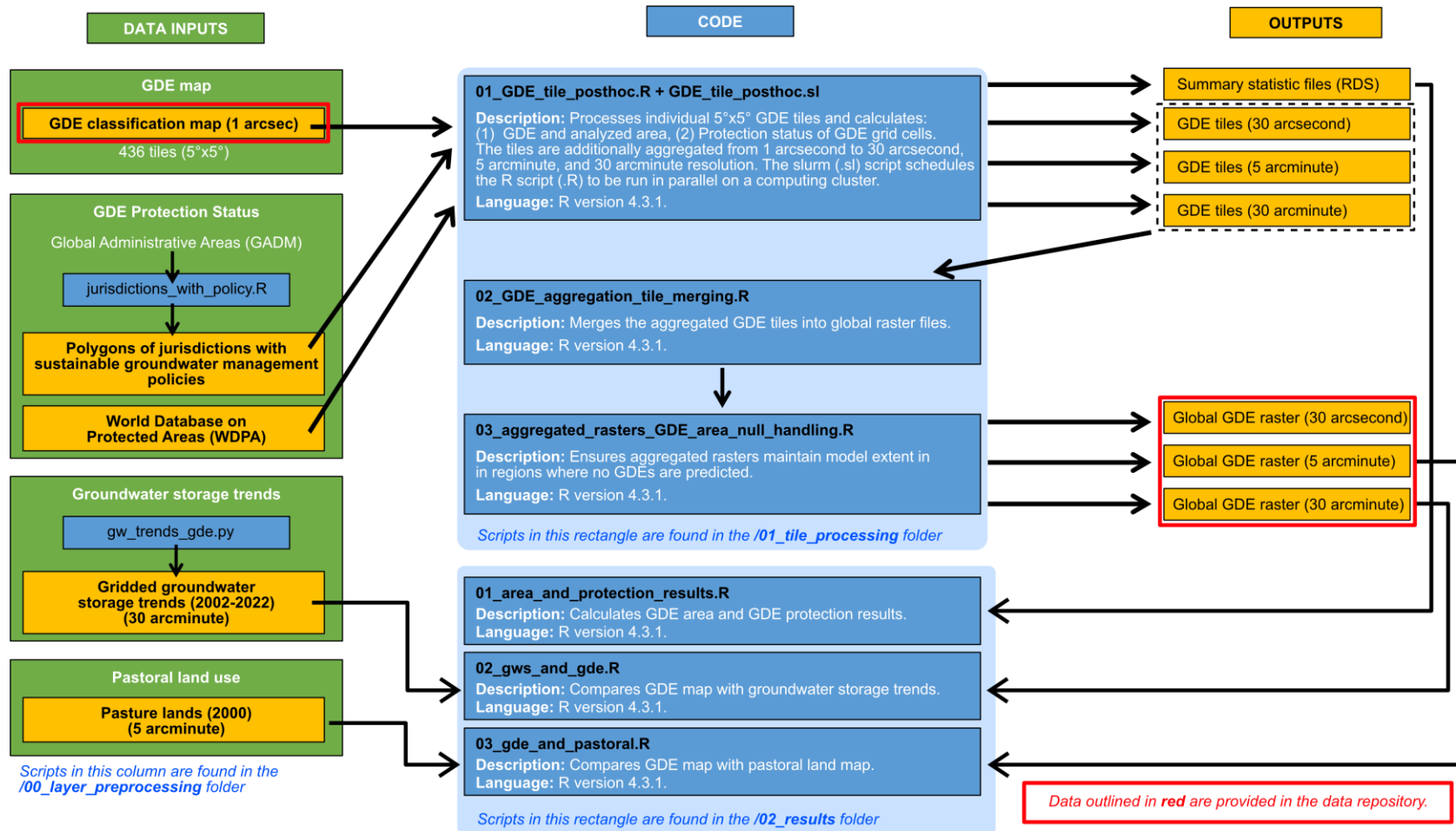
**B9 - Extract predictor variable values for each random global point from model extent**

<b>Script name</b>	B9_GlobalRandomPointSampler.js
<b>Description</b>	This Google Earth Engine code extracts predictor values at randomly generated points within the model extent to compare distributions with those of training points in the model.
<b>Data inputs</b>	All data inputs are uploaded and accessible as Google Earth Engine assets. Links are provided in the code.  Masking Data: <ul style="list-style-type: none"> <li>• GlobalDTG</li> <li>• ESRI 10 m Land Use Land Cover data</li> <li>• Dryland Regions (Beck_KG_V1_present_0p0083)</li> </ul> Predictor Variables: <ul style="list-style-type: none"> <li>• Compound Topographic Index (TopConvIndex_Global_KGclims)</li> <li>• Climate - precipitation and potential evapotranspiration (IDAHO_EPSCOR/TERRACLIMATE)</li> <li>• Actual Evapotranspiration (CAS/IGSNRR/PML/V2_v017)</li> <li>• Landsat imagery (LANDSAT/LC08/C02/T1_L2)</li> <li>• randomGlobalPts_validationN</li> <li>• randomGlobalPts_validationS</li> </ul>
<b>Intermediary data outputs</b>	globalvalpts1_withpredictors.csv globalvalpts2_withpredictors.csv globalvalpts3_withpredictors.csv globalvalpts4_withpredictors.csv globalvalpts5_withpredictors.csv globalvalPts.csv
<b>Authors</b>	Melissa M. Rohde

**B10 - Generate distribution plots of predictor values for the global and model training points.**

<b>Script name</b>	B10_Predictor_DistributionPlots.Rmd
<b>Description</b>	Generates overlapping distribution plots and statistic.
<b>Data inputs</b>	globalvalpts1_withpredictors globalvalpts2_withpredictors globalvalpts3_withpredictors globalvalpts4_withpredictors globalvalpts5_withpredictors globalvalPts  trainingpts1_withpredictors trainingpts2_withpredictors trainingpts3_withpredictors trainingpts4_withpredictors trainingpts5_withpredictors trainingpts6_withpredictors trainingpts7_withpredictors trainingpts8_withpredictors trainingPts  validationpts1_withpredictors validationpts2_withpredictors  joinval.csv jointrain.csv
<b>R packages used</b>	<i>overlapping</i> (Pastore et al., 2022). <i>gridExtra</i> (Auguie & Antonov, 2017). <i>tidyverse</i> (Wickham et al., 2019)
<b>Intermediary data outputs</b>	Distribution plots and statistics (Extended Data Figure 5).
<b>Authors</b>	Melissa M. Rohde

## Section C Post hoc analysis and result reporting



**Section C**  
**Post hoc analysis and result reporting**

SUBSECTION:  
**POSTHOC LAYER PREPROCESSING**

<b>Script name</b>	00_layer_preprocessing/gw_trends_gde.py
<b>Description</b>	Derives GRACE-based gridded groundwater storage trends at 30 arcminute resolution using GLDAS-2.1 output from both Noah and VIC land surface models. The script provided shows the workflow using Noah model output, however the same is also performed with VIC model output.
<b>Python packages used</b>	<i>NumPy</i> (Harris et al., 2020). <i>xarray</i> (Hoyer & Harmann, 2017). <i>mvstats</i> (Chandanpurkar, 2018).
<b>Intermediary data outputs</b>	gldas_2.1_vic_local_trends.nc gldas_2.1_noah_local_trends.nc
<b>Authors</b>	Hrishikesh A. Chandanpurkar

<b>Script name</b>	00_layer_preprocessing/jurisdictions_with_policy.R
<b>Description</b>	Generates a vector file representing jurisdictions with sustainable groundwater management policies.
<b>Python packages used</b>	<i>terra</i> (Hijmans, 2023).
<b>Intermediary data outputs</b>	all_juris_with_GDE_protection_policies.sqlite
<b>Authors</b>	Xander Huggins

**Section C:**  
**Post hoc analysis and result reporting**

SUBSECTION:  
**TILE PROCESSING**

<b>Script name</b>	01_tile_processing/01_GDE_tile_posthoc.R
<b>Description</b>	For each 5-degree GDE tile at 1 arcsecond resolution, this script: <ol style="list-style-type: none"> <li>1) Calculates grid cell area.</li> <li>2) Computes GDE area, analyzed area, and grid cell area.</li> <li>3) Rasterizes the WDPA and extent of jurisdictions with sustainable groundwater management policies at 1 arcsecond resolution.</li> <li>4) Zonally summarizes GDE area within different forms of protection.</li> <li>5) Computes GDE area and analyzed area per continent.</li> <li>6) Aggregates GDE area, analyzed area, and grid cell area to the resolutions of: 30 arcseconds, 5 arcminutes, and 30 arcminutes.</li> <li>7) Computes (i) GDE area to analyzed area, (ii) GDE area to grid cell area, and (iii) analyzed area to grid cell area fractions.</li> </ol>
<b>R packages used</b>	<i>terra</i> (Hijmans, 2023). <i>raster</i> (Hijmans, 2023). <i>rasterDT</i> (Hijmans, 2023). <i>readr</i> (Wickham et al., 2023).
<b>Intermediary data outputs</b>	Zonal summary statistic file per tile. Aggregated GDE tiles at 30 arcsecond, 5 arcminute, and 30 arcminute resolution.
<b>Authors</b>	Xander Huggins

<b>Script name</b>	01_tile_processing/GDE_tile_posthoc.sl
<b>Description</b>	Slurm script to run GDE tile post hoc analyses on cluster. The R script called in this slurm script is provided below.
<b>Cluster acknowledgement</b>	Computing resources were provided by the Digital Research Alliance of Canada ( <a href="https://alliancecan.ca/">https://alliancecan.ca/</a> ).
<b>Authors</b>	Xander Huggins

<b>Script name</b>	01_tile_processing/ <b>02_GDE_aggregation_tile_merging.R</b>
<b>Description</b>	Merges the exported aggregated (i.e., 30s, 5m, 30m) tiles into global (single) rasters.
<b>R Packages</b>	<i>terra</i> (Hijmans, 2023).
<b>Authors</b>	Xander Huggins

<b>Script name</b>	01_tile_processing/ <b>03_aggregated_rasters_GDE_area_null_handling.R</b>
<b>Description</b>	Ensures that grid cells with GDE area = 0 and analyzed area > 0 are reported with GDE areas and GDE area fractions of 0 rather than NA.
<b>Data outputs</b>	Global aggregated GDE rasters at 30 arcsecond, 5 arcminute, and 30 arcminute resolution for deposition on the data repository.
<b>R Packages</b>	<i>terra</i> (Hijmans, 2023).
<b>Authors</b>	Xander Huggins

**Section C:**  
**Post hoc analysis and result reporting**

SUBSECTION:  
**RESULT REPORTING**

<b>Script name</b>	02_results/ <b>01_area_and_protection_results.R</b>
<b>Description</b>	Reports base GDE results such as GDE area and analyzed area per continent, and GDE protection stats.
<b>R Packages</b>	<i>terra</i> (Hijmans, 2023).
<b>Authors</b>	Xander Huggins

<b>Script name</b>	02_results/ <b>02_gws_and_gde.R</b>
<b>Description</b>	Compares GDE area density to groundwater storage trends at 30 arcminute resolution. Develops the bi-variate map as plotted in Figure 2. Calculates area-weighted groundwater storage trends and GDE area density for select freshwater ecoregions of the world.
<b>R Packages</b>	<i>terra</i> (Hijmans, 2023). <i>tidyverse</i> (Wickham et al., 2019).
<b>Authors</b>	Xander Huggins

<b>Script name</b>	02_results/ <b>03_gde_and_pastoral.R</b>
<b>Description</b>	Develops bi-variate map as plotted in Extended Data Figure 7, comparing GDE area density with pasture land area density. Calculates percentage of GDE area that lies within regions with at least 25% pasture land area density.
<b>R Packages</b>	<i>terra</i> (Hijmans, 2023). <i>tidyverse</i> (Wickham et al., 2019).
<b>Authors</b>	Xander Huggins



## Section D

### Data deposit and access

Global GDE data are deposited and openly accessible at four resolutions. All raster layers are referenced to the World Geodetic System 1984 (WGS84) coordinate reference system.

#### Core GDE data:

##### Dataset 1: 1 arcsecond GDE data (~30 m grids at the equator)

At this base resolution, two raster layers are provided:

Layer: **classification**

Values:

0: Pixel is outside of model domain

1: Pixel is likely a GDE

2: Pixel is not likely a GDE

Layer: **probability**

Values:

[0-100]: Likelihood pixel is a GDE (100) or non-GDE (0)

These layers are provided across 436 tiles, each with a 5° x 5° extent. For example, the tile named “n45w105.tif” has an extent of: xmin: -105, xmax: -100, ymin: 45, ymax: 50.

Access: These layers can be accessed individually on the repository by navigating to: [./GDE\\_tiles/](#)

#### Aggregated GDE data:

##### Dataset 2: Aggregated GDE data at 30 arcsecond resolution (~1 km grids at the equator)

##### Dataset 3: Aggregated GDE data at 5 arcminute resolution (~10 km grids at the equator)

##### Dataset 4: Aggregated GDE data at 30 arcminute resolution (~50 km grids at the equator)

At these aggregated resolutions, five raster layers are provided. To manage file size of global rasters, all layers are provided in the INT4U datatype (i.e., capable of storing integers from 0 to  $2^{32}$ ). Operating within this upper constraint, fraction layers (normally ranging [0-1]) are multiplied by  $10^8$  to retain the greatest precision possible. Area layers, across all resolutions, do not have grid cell values that exceed the upper limit of the datatype and thus are not modified other than providing area values that are rounded to the nearest square meter. All aggregated grid cells with no analyzed area are set to NA.

Layer: **GDE\_sqm**

Values:

GDE area (m<sup>2</sup>) within grid cell.

Layer: **GDE\_frac\_AA**

Values:

[0- $10^8$ ]: Fraction of analyzed area (i.e., area within the model domain) that is a GDE. This fraction is multiplied by  $10^8$  and saved as an integer. To convert this layer to the range [0-1], divide by  $10^8$ .

Layer: **GDE\_frac\_GA**

Values:

[0-10<sup>8</sup>]: Fraction of grid cell area that is a GDE. This fraction is multiplied by 10<sup>8</sup> and saved as an integer. To convert this layer to the range [0-1], divide by 10<sup>8</sup>.

Layer: **AA\_sqm**

Values:

Analyzed area (m<sup>2</sup>) within grid cell, rounded to nearest square meter.

Layer: **AA\_frac\_GA**

Values:

[0-1e8]: Fraction of grid cell area that is analyzed (i.e., within the model domain). This fraction is multiplied by 10<sup>8</sup> and saved as an integer. To convert this layer to the range [0-1], divide by 10<sup>8</sup>.

Access: These layers can be accessed as consistent global rasters (i.e., individual tiles have been mosaiced) by navigating in the repository to: `./GDE_aggregated_layers/`

## License agreement

By using any of these datasets, you agree to cite Rohde et al. (2024) – see below – in publications that make use of any of the above datasets. These data are licensed under Creative Commons Attribution 4.0 International (CC BY 4.0) license. To view a copy of this license, visit:

<https://creativecommons.org/licenses/by/4.0/>

## Data citation

Rohde, M.M., C.M. Albano, X. Huggins, K.R. Klausmeyer, C. Morton, A. Sharman, E. Zaveri, L. Saito, Z. Freed, J.K. Howard, N. Job, H. Richter, K. Toderich, A. Rodella, T. Gleeson, J. Huntington, H.A. Chandanpurkar, A.J. Purdy, J.S. Famiglietti, M.B. Singer, D.A. Roberts, K. Caylor, J.C. Stella. Groundwater-dependent ecosystem map exposes global dryland protection needs. *Nature* (2024). DOI: [10.1038/s41586-024-07702-8](https://doi.org/10.1038/s41586-024-07702-8)

## References

- Auguie, B., & Antonov, A. gridExtra: Miscellaneous Functions for "Grid" Graphics. R package version 2.3. <https://cran.r-project.org/package=gridExtra> (2017).
- Chandanpurkar, H.A. mvstats: Vectorized multivariate statistical functions for analyzing multi-dimensional earth system data. Python package version 0.1.0. <https://mvstats.readthedocs.io> (2018).
- Cohen, A. fuzzywuzzy: Fuzzy string matching in python. Python package version 0.18.0. <https://pypi.org/project/fuzzywuzzy> (2020).
- Fan, Y., Miguez-Macho, G., Jobbágy, E. G., Jackson, R. B. & Otero-Casal, C. Hydrologic regulation of plant rooting depth. *Proc National Acad Sci* **114**, 10572–10577 (2017).
- Harris, C.R., Millman, K.J., van der Walt, S.J. et al. Array programming with NumPy. *Nature* **585**, 357–362 (2020).
- Hijmans, R.J.. raster: Geographic Data Analysis and Modeling. R package version 3.6-20, <https://cran.r-project.org/package=raster> (2023)
- Hijmans, R. J. terra: Spatial Data Analysis. R package version 1.7-37 <https://CRAN.R-project.org/package=terra> (2023).
- Hoyer, S. & Hamman, J. xarray: N-D labeled Arrays and Datasets in Python. *Journal of Open Research Software* **5**(1), 1-10. (2017).
- Klausmeyer K., J. Howard, T. Keeler-Wolf, K. Davis-Fadtke, R. Hull, A. Lyons. Mapping Indicators of Groundwater Dependent Ecosystems in California: Methods Report. San Francisco, California. (2018).
- McKinney, W. Data structures for statistical computing in python. *Proceedings of the 9th Python in Science Conference*, 56-61 2010).
- O'Brien, J. rasterDT: Fast Raster Summary and Manipulation. R package version 0.3.2. <https://CRAN.R-project.org/package=rasterDT> (2022).
- Pandas development team. *pandas*. Python package version 1.5.2. <https://pandas.pydata.org/> (2020).
- Pastore, M., Di Loro, P.A., Mingione, M., Calcagni, A. overlapping: Estimation of overlapping empirical distributions. R package version 2.1. <https://cran.r-project.org/package=overlapping> (2022).
- Sabatini, F.M., Lenoir, J., Bruehlheide, H. & the sPlot Consortium. sPlotOpen – An environmentally-balanced, open-access, global dataset of vegetation plots. [Dataset] <https://doi.org/10.25829/ivid.3474-40-3292> (2021).
- Wickham H., Averick M., Bryan J., et al. "Welcome to the tidyverse." *Journal of Open Source Software*, **4**(43), 1686 (2019).
- Wickham, H., François, R., Henry, L., et al. dplyr: A Grammar of Data Manipulation. R package version 1.1.2 <https://cran.r-project.org/package=dplyr> (2023).
- Wickham, H., Hester, J., Bryan, J. readr: Read Rectangular Text Data. R package version 2.1.4, <https://cran.r-project.org/package=readr> (2023).