



1 Spatially explicit re-harmonized terrestrial carbon densities for 2 calibrating Integrated human-Earth System Models

3 Authors- Kanishka B. Narayan^{i*}, Alan V. Di Vittorioⁱⁱ, Evan Margiottaⁱ, Seth A. Spawn-Lee^{iii, iv}, Holly K. Gibbs^{iii, iv}

4 (i) Joint Global Change Research Institute (JGCRI), Pacific Northwest National Lab (PNNL), Washington DC, USA

5 (ii) Lawrence Berkeley National Lab (LBNL), Berkeley, CA, USA

6 (iii) Department of Geography, University of Wisconsin-Madison, Madison, WI, USA

7 (iv) Center for Sustainability and the Global Environment (SAGE), Nelson Institute for Environmental Studies, University of
8 Wisconsin-Madison, Madison, WI, USA

9 *Correspondence to: Kanishka B. Narayan (kanishka.narayan@pnnl.gov)

10 **Abstract-** Soil and vegetation carbon densities play a critical role in global and regional human-
11 Earth system models. These densities affect variables such as land use change emissions and also
12 influence land use change pathways under climate forcing scenarios where terrestrial carbon is
13 assigned a carbon price. Recently, more spatially explicit, fine resolution data have become
14 available for both soil and vegetation carbon. However, for models to effectively use these data
15 the fine resolution data need to be reharmonized to the initial land use and land cover conditions
16 represented by these models. Without such reharmonization the carbon values may be inaccurate
17 for particular land types and places where the source data and the model disagree on the land
18 use/cover type. Here we present reharmonized soil and vegetation carbon densities both at the 5-
19 arcmin resolution grid cell level and also aggregated to 235 water sheds for 4 land use types and
20 15 land cover types. These data are particularly useful as initial land carbon conditions for global
21 Multisectoral Dynamic Models (MSD). Moreover, these data include six different statistical
22 states calculated using distinct resampling methods for each of the land use and land cover types.
23 These statistical states are used to define a range of possible carbon values for each land
24 classification, and any state can be used for defining initial conditions of soil and vegetation
25 carbon in MSD models. Users can also estimate any percentile of the carbon distribution defined
26 by these six summary states. We make use of these statistical states to calculate spatially distinct
27 uncertainties in the carbon densities by land type. We have implemented these data in a state-of-
28 the-art multi sector dynamics model, namely the Global Change Analysis Model (GCAM), and
29 show that these new data improve several land use responses in the model, especially when
30 terrestrial carbon is assigned a carbon price. The statistical states in our data are validated against
31 similar estimates in the literature both at a grid cell level and at a regional level.

32 1. Introduction

33 Soil and vegetation carbon densities play a critical role in global and regional models such as
34 Earth system models (ESMs), multisector dynamics models (MSDs) and integrated human-Earth
35 system models. These densities influence the predicted productivity of land types (e.g., forest
36 yields, pasture yields, and crop yields) and directly influence land use change emissions.
37 Moreover, these densities affect land use change pathways under climate forcing scenarios
38 implemented in these models (Thomson et al., 2010; Wise et al., 2009). Many models make use
39 of carbon density data that are differentiated by land type but are not spatially explicit. For
40 example, models have previously used estimates of carbon values on undisturbed land from
41 Houghton et al. (Houghton, 1999) and the IPCC (Jackson et al., 2017), among others. Recently,



1 spatially explicit soil carbon density data have been made available by the FAO (Nachtergaele et
2 al., 2010) at a 1 km resolution and at 250 m resolution by the SoilGrids team at the International
3 Soil Reference and Information Centre (ISRIC) (Batjes et al., 2017; Hengl et al., 2014).
4 Similarly, spatially explicit data on vegetation carbon differentiated by above and below ground
5 biomass and spanning several vegetation types have been made available by Spawn et al. (Spawn
6 et al., 2020). Use of spatially distinct, fine resolution data has the potential to significantly
7 improve results from global and regional models by better capturing the geographies of soil and
8 vegetation carbon stocks (Jungkunst et al., 2022). These data can also be used to explore and
9 validate effects of different carbon parameterizations in models (Wieder et al., 2014).

10 However, these carbon data need to be transformed significantly in order to be used in a robust
11 manner by regional and global models. This is because each of these fine resolution datasets
12 utilizes its own assumptions of land use and land cover which may be distinct from the land use
13 and land cover definitions used by the models in question. For example, many of these fine
14 resolution data use land cover definitions from the European Space Agency Climate Change
15 Initiative (ESA CCI) dataset (Li et al., 2018; Liu et al., 2018) while models may use land use
16 definitions from the Historical Database of the Global Environment (HYDE) dataset (Klein
17 Goldewijk et al., 2017) and/or land cover definitions from the Moderate Resolution Imaging
18 Spectrometer (MODIS) (Barnes et al., 2003; Justice et al., 2002).

19 Resolution mismatch between data and models provides an additional challenge. The new,
20 spatially distinct carbon densities are available at a very fine resolution (250m / 300 m) while
21 models are often configured to use coarser data that better match their working resolution. For
22 example, consistent land datasets that have frequently been used for climate modelling are
23 available at a resolution of 5 arcmins (i.e. ~ 10km at the equator) (van Asselen & Verburg,
24 2012), and many regional models operate on land units defined by geopolitical and/or
25 geophysical boundaries. Given the difference in resolutions, and the above-mentioned
26 differences in land classifications, a consistent harmonization method is required to appropriately
27 match the fine resolution carbon data with the appropriate land uses and land cover types within
28 a model.

29 Here, we prepare and present an aggregated dataset of fine resolution carbon density for soil and
30 vegetation biomass in MgC/ha that has been aligned with the land use and land cover definitions
31 and distribution in the Global Change Analysis Model (GCAM) (Calvin et al., 2019). GCAM
32 represents the interactions between five major systems – energy, water, land, climate, and the
33 economy – at global and regional scales. The soil data are based on the 250 m-resolution
34 SoilGrids dataset and represent a depth of 0-30 cms (Hengl et al., 2014). The aboveground and
35 below ground biomass are based on the 300 m-resolution Spawn et al. dataset (Spawn et al.,
36 2020). The land data integration for GCAM is performed at a resolution of 5 arcmin using a grid-
37 based land data system (Di Vittorio et al., 2020). Hence the final outputs are available as rasters
38 at 5 arcmin resolution. These outputs include rasters corresponding to each of the land use and
39 land cover types in GCAM.

40 In addition to the reharmonization, we calculated six different data driven statistical states for
41 each carbon pool and each 5 arcmin grid cell using different resampling methods (area weighted



1 average, minimum, maximum, median, Q1, Q3) when re-gridding the data. We also present an
2 easy-to-use tabular output summarizing the six carbon density states for each carbon pool within
3 each of the 235 watersheds intersected with 207 country (ISO) boundaries that are modelled by
4 GCAM. Using our six statistical states, users can calculate any percentile within the distribution
5 of carbon both at a pixel level and at a land region/ water basin level for any land type (For
6 example, the 90th percentile). Such a calculation would not be time intensive given that the six
7 summary states are already available at multiple scales.

8 GCAM needs to be initialized with densities that represent long term potential maximum carbon
9 values since these values are used to spin-up the model in historical years. In particular, the
10 density values are used to spin up the carbon cycle from 1700-1975. Various studies have found
11 that the long term potential carbon densities are much higher than the contemporary values due
12 to ongoing land use and cover change (Erb et al., 2018; Walker et al., 2022). Moreover, studies
13 have also highlighted difficulties in estimating long term potential carbon densities, since these
14 estimations require long spin up periods themselves (Fang et al., 2014). Here we addressed this
15 issue by deriving a more data driven long term potential carbon state from our new dataset. By
16 analyzing the distribution of carbon values within each land type-watershed combination we
17 found two potential options for initializing GCAM. The first is the Q3 state which would
18 represent a low carbon initialization in 1700 (This state results in 2144 PgC of terrestrial carbon
19 in 1700) and the second, the 90th percentile state which would represent a high carbon
20 initialization (which results in an initial terrestrial carbon stock of 3028 PgC in 1700). We find
21 that utilizing this new carbon dataset for the spin-up improved several responses in GCAM,
22 especially under forcing scenarios where the value of terrestrial carbon is priced using a carbon
23 tax.

24 We also compared the Q3 and the 90th percentile carbon state in our dataset (which are intended
25 to represent a pre-industrial carbon state) with estimates of potential pre-industrial top-soil
26 carbon by grid cell from Sanderman et al. (Sanderman et al., 2017) and with similar estimates of
27 vegetation carbon from Walker et al. (Walker et al., 2022). We also perform global-level
28 validation of our carbon data, respecting that there is a high degree of uncertainty in carbon
29 estimates from different datasets (Scharlemann et al., 2014; Tifafi et al., 2018). We compare the
30 global estimates of carbon from our reharmonized data with similar estimates of soil and
31 vegetation carbon from other sources in the literature, with different meta-analyses of carbon
32 inventories (Scharlemann et al., 2014) and with modelled estimates of contemporary and
33 historical soil carbon (Sanderman et al., 2017).

34 The available dataset includes raster files for the six different statistical states for each land use
35 type (Cropland, Urban land, Pasture and Unmanaged land) and each carbon pool, bringing the
36 total to 72 distinct raster files. We also provide a thematic file that labels each cell with the
37 dominant biome for Unmanaged land (out of 15). We also present a tabulated text file with the
38 six carbon state values for each land type and carbon pool aggregated to 699 land regions (235
39 water basins intersected with 207 country boundaries). Making the data available at these
40 different resolutions should help facilitate effective multiscale modelling of terrestrial carbon.
41 We implemented this carbon reharmonization programmatically in a land data system, *moirai*



1 (Di Vittorio et al., 2020), which can further be used to update the data, validate the data, and
 2 reharmonize the data to any other land use and land cover types required by models other than
 3 GCAM.

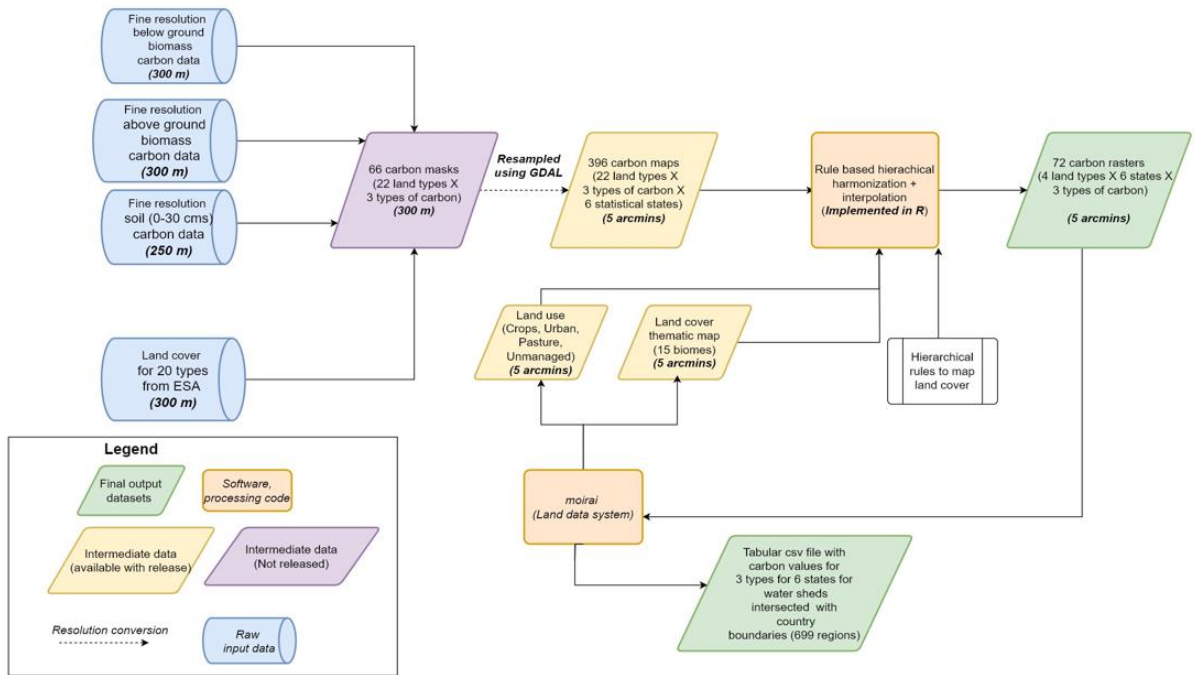
4

5 **2. Description of data processing**

6 Our carbon data processing method can be organized into three stages:

- 7 1. Stage 1- Resampling source datasets based on fine resolution land cover
- 8 2. Stage 2- Re-mapping the carbon to Moirai land use and land cover
- 9 3. Stage 3- Aggregating raster carbon data to 699 land regions

10 Figure 1 below summarizes our processing approach from start to finish



11

12 *Figure 1: Description of data processing implementation to generate carbon datasets*

13

14 **2.1 Stage 1 – resampling source data**

15 This stage combines the soil and vegetation carbon data (Mg/ha) both at 300 m resolution with
 16 the input land cover assumptions from the ESA CCI dataset that correspond to these data. Note
 17 that since the land cover dataset from ESA CCI is at a 300m resolution, we resample both our
 18 carbon datasets to 300m before this stage.

19



1 We first generate land cover masks (1=respective land type present, 0=otherwise) for each of 22
 2 aggregated ESA CCI land cover types (**Table 1**). We combine the land cover masks with the
 3 carbon data to create 66 rasters (22 land types X 3 carbon pools), each representing a carbon data
 4 mask for an ESA land type. The resulting rasters are calculated as follows:

$$5 \quad Carbon_LT_300m_{pool,j,LT} = Carbon_300m_{pool,j} * LT_mask_300m_j \quad (1)$$

6 Where,

7

8 j is the index of a 300m grid cell,

9 pool is the carbon pool (soil, aboveground biomass, belowground biomass),

10 LT is the ESA land type.

11

12 Next we use six distinct resampling methods to re-grid these data to a 5 arcmin resolution. Each
 13 method is applied to each of the land types and thus we derive 6 statistical states for each 5
 14 arcmin grid cell. These aggregated rasters are calculated as follows:

15

$$16 \quad Carbon_LT_5arcmin_{pool,i,state} =$$

$$17 \quad state \left(\begin{matrix} Carbon_LT_300m_{pool,j} & Carbon_LT_300m_{pool,j+2} \\ Carbon_LT_300m_{pool,j+1} & Carbon_LT_300m_{pool,j+n} \end{matrix} \right) \quad (2)$$

18

19 Where,

20

21 i is the index of a 5 arcmin grid cell,

22 pool is the carbon pool (soil, aboveground biomass, belowground biomass),

23 state is the resampling method (weighted average, median, min, max, q1, q3),

24 j is the index of each 300 m grid cell within aggregated cell i,

25 n is the total number of 300 m cells that are aggregated into cell i.

26

27 Thus, we generate 366 (22 land cover types X 3 types of carbon X 6 states) layers of carbon that
 28 correspond to the aggregated ESA CCI land cover types. This processing is largely conducted
 29 through the GDAL software (Warmerdam, 2008) and implemented using bash scripts.

30

31

ESA_entry code	ESA_classes
10;Cropland_rainfed;255;255;100	Cropland
11;Herbaceous cover;255;255;100	Unknown_Herb
12;Tree or shrub cover;255;255;0	Unknown_Tree
20;Cropland_irrigated or post-flooding;170;240;240	Cropland



30;Mosaic cropland (>50%) / natural vegetation (tree shrub herbaceous cover) (<50%);220;240;100	Mosaic_Crop
40;Mosaic natural vegetation (tree shrub herbaceous cover) (>50%) / cropland (<50%) ;200;200;100	Mosaic_tree
50;Tree cover broadleaved evergreen	Broadleaf_Evergreen
60;Tree cover broadleaved deciduous	Broadleaf_Deciduous
61;Tree cover broadleaved deciduous	Broadleaf_Deciduous
62;Tree cover broadleaved deciduous	Broadleaf_Deciduous
70;Tree cover needleleaved evergreen	Needleleaved_Evergreen
71;Tree cover needleleaved evergreen	Needleleaved_Evergreen
72;Tree cover needleleaved evergreen	Needleleaved_Evergreen
80;Tree cover needleleaved deciduous	Needleleaved_deciduous
81;Tree cover needleleaved deciduous	Needleleaved_deciduous
82;Tree cover needleleaved deciduous	Needleleaved_deciduous
90;Tree cover mixed leaf type (broadleaved and needleleaved);120;130;0	Mixed_Forests
100;Mosaic tree and shrub (>50%) / herbaceous cover (<50%);140;160;0	Mosaic_tree
110;Mosaic herbaceous cover (>50%) / tree and shrub (<50%);190;150;0	Mosaic_Herb
120;Shrubland;150;100;0	Shrubland
121;Shrubland evergreen;120;75;0	Shrubland
122;Shrubland deciduous;150;100;0	Shrubland
130;Grassland;255;180;50	Grasslands
140;Lichens and mosses;255;220;210	Grasslands
150;Sparse vegetation (tree shrub herbaceous cover) (<15%);255;235;175	Sparse_Tree
151;Sparse tree (<15%);255;200;100	Sparse_Tree
152;Sparse shrub (<15%);255;210;120	Sparse_Shrub
153;Sparse herbaceous cover (<15%);255;235;175	Sparse_Shrub
160;Tree cover flooded fresh or brakish water;0;120;90	Flood_Tree_Cover
170;Tree cover flooded saline water;0;150;120	Flood_Tree_Cover
180;Shrub or herbaceous cover flooded fresh/saline/brakish water;0;220;130	Flooded_Shrub
190;Urban areas;195;20;0	Urbanland
200;Bare areas;255;245;215	Desert
201;Consolidated bare areas;220;220;220	Desert
202;Unconsolidated bare areas;255;245;215	Desert
220;Permanent snow and ice;255;255;255	Polar_desert_rock_ice

1
 2
 3

Table 1: Raw ESA codes mapped to ESA land types



1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16

2.2 Stage 2 – remapping the carbon data to Moirai land use/cover

2.2.1 Reharmonization of ESA land cover with Moirai land cover at 5 arcmins using a prioritization matrix

Next, the 366 layers described above are aligned with the default initial land use/cover for GCAM (2010) at a 5 arcmin resolution. These initial land use/cover data are generated by the Moirai land data system based on land use data from the HYDE (Klein Goldewijk et al., 2017) database and a one-half degree land cover product (Meiyappan & Jain, 2012). Moirai can generate land use and land cover maps for any year based on the these datasets combined with a potential vegetation dataset from Ramankutty et al. (1999). The Moirai land use and land cover types are listed in table 2. It is important to note that carbon values are independently assigned to each of the four Moirai land use types in each cell, and that the unmanaged land use type can be only one of the Moirai land cover types in a given cell. Moirai is described in more detail in Di Vittorio et al. (Di Vittorio et al., 2020).

Land use	Land cover
Cropland	Cropland
Pasture	Pasture
Urbanland	Urbanland
Unmanaged	TropicalEvergreenForest/Woodland
	TropicalDeciduousForest/Woodland
	TemperateBroadleafEvergreenForest/Woodland
	TemperateNeedleleafEvergreenForest/Woodland
	TemperateDeciduousForest/Woodland
	BorealEvergreenForest/Woodland
	BorealDeciduousForest/Woodland
	Evergreen/DeciduousMixedForest/Woodland
	Savanna
	Grassland/Steppe
	DenseShrubland
	OpenShrubland
	Tundra
Desert	
PolarDesert/Rock/Ice	

17
18
19
20
21
22

Table 2: land use, land cover types for Moirai/GCAM. Total of 4 land use types, 15 types of land cover tracked for Unmanaged land type

The carbon for each Moirai land type in a cell needs to be selected from the 366 rasters generated in Stage 1 described above. We use a rule-based harmonization approach where we select the



1 appropriate carbon values by matching the Moirai land type with the corresponding ESA land
2 cover type (Table 3). We assign 6 possible ESA land cover types to each Moirai land type and
3 rank them according to their similarity with the Moirai land type. This means that carbon values
4 for a particular moirai land type can come from any of six ESA land cover types, as long as they
5 are present in a given cell. For example, a Tropical Evergreen Forest cell in Moirai, may be
6 assigned carbon values from the Evergreen_Combined, Mixed_Forests, Mosaic_Tree,
7 Flood_Tree_Cover, Unknown_Tree_Cover, or Sparse_Treecover ESA land cover types. The
8 similarity ranking both maximizes the number of Moirai land type assignments and ensures that
9 the most appropriate carbon values are selected. The first ESA land cover in the ranked list that is
10 present in a given cell provides the carbon values for the corresponding Moirai land type in the
11 same cell (Table 3). In the example above, The Evergreen_Combined carbon data would be
12 chosen first over all other ESA land covers if it existed in a given cell and the Sparse_Treecover
13 carbon data would be chosen if it were the only ESA land cover from the list that existed in a
14 given cell. These prioritization rules are designed such that carbon data from one biome is not
15 assigned to a different biome when reharmonizing and re-gridding the carbon. The ESA land
16 cover selection is done once for each cell and Moirai land type, and then the data from the
17 corresponding carbon pool and state rasters are assigned to the Moirai land type in the target cell.
18 This results in 72 rasters that become input files for Moirai.

19
20
21

22 We used expert judgement when developing the matrix so as to best represent the Moirai land
23 types when selecting from the ESA land types. For certain land types we constrain the choices by
24 allowing less than six choices. For example, carbon for a Moirai Desert cell can only be chosen
25 from a corresponding desert cell in the ESA masks. On the other hand, Moirai Tundra includes
26 eight ESA land covers because ESA does not have an explicit Tundra class. The increased
27 number of options aims to provide adequate data coverage for Tundra¹. Furthermore, certain
28 biome types that are not modelled by GCAM or represented explicitly in Moirai receive low
29 priority rankings. For example, Flooded land types are never included as a first priority choice for
30 any land type since Moirai does not explicitly include flooded land types. Conversely, the ESA
31 land cover data do not include any explicit representation of pastures or rangeland. Our rules
32 assign pasture carbon values based on proximate grassland or shrubland carbon values.
33 Grasslands in particular are prioritized for Pasture carbon selection because the pasture definition
34 in GCAM corresponds to grasslands used for grazing.

35

¹ Tundra data selection prioritizes polar desert rock ice pixels. The location of these pixels coincides with the Tundra land cover and they also represent pixels with high values for soil carbon densities.



moirai land type name	Corresponding ESA LAND COVER prioritized from Primary to 8th							
	Primary	2	3	4	5	6	7	8
TropicalEvergreenForest/Woodland	Evergreen_Combined	Mixed_Forests	Mosaic_Tree	Flood_Tree_Cover	Unknown_Tree	Sparse_Tree	-	-
TropicalDeciduousForest/Woodland	Deciduous_Combined	Mixed_Forests	Mosaic_Tree	Flood_Tree_Cover	Unknown_Tree	Sparse_Tree	-	-
TemperateBroadleafEvergreenForest/Woodland	Broadleaf_Evergreen	Mixed_Forests	Mosaic_Tree	Flood_Tree_Cover	Unknown_Tree	Sparse_Tree	-	-
TemperateNeedleleafEvergreenForest/Woodland	Needleleaved_Evergreen	Mixed_Forests	Mosaic_Tree	Flood_Tree_Cover	Unknown_Tree	Sparse_Tree	-	-
TemperateDeciduousForest/Woodland	Deciduous_Combined	Mixed_Forests	Mosaic_Tree	Flood_Tree_Cover	Unknown_Tree	Sparse_Tree	-	-
BorealEvergreenForest/Woodland	Evergreen_Combined	Mixed_Forests	Mosaic_Tree	Flood_Tree_Cover	Unknown_Tree	Sparse_Tree	-	-
BorealDeciduousForest/Woodland	Combined_Deciduous	Mixed_Forests	Mosaic_Tree	Flood_Tree_Cover	Unknown_Tree	Sparse_Tree	-	-
Evergreen/DeciduousMixedForest/Woodland	Mixed_Forests	Mosaic_Tree	Flood_Tree_Cover	Unknown_Tree	Sparse_Tree	Sparse_Tree	-	-
Savanna	Mosaic_Herb	Grasslands	Unknown_Herb	Flood_Shrub	-	-	-	-
Grassland/Steppe	Grasslands	Unknown_Herb	Mosaic_Herb	Flood_Shrub	-	-	-	-
DenseShrubland	Shrubland	Unknown_Tree	Flooded_Shrub	Mosaic_Herb	-	-	-	-
OpenShrubland	Sparse_Shrub	Mosaic_Herb	Flooded_Shrub	Unknown_Herb	-	-	-	-
Tundra	Polar_Desert_Rock_Ice	Sparse_Shrub	Mosaic_Herb	Unknown_Herb	Unknown_Tree	Sparse_Tree	Shrubland	Mosaic_Tree
Desert	Desert	-	-	-	-	-	-	-
PolarDesert/Rock/Ice	Polar_desert_rock_ice	-	-	-	-	-	-	-
Cropland	Cropland	Mosaic_Cropland	-	-	-	-	-	-
Pasture	Grasslands	Mosaic_Herb	Unknown_Herb	Sparse_Tree	Sparse_Shrub	-	-	-
Urbanland	Urbanland	-	-	-	-	-	-	-

1

2

Table 3: Prioritization matrix to match ESA land cover with moirai land types

3

4

5

2.2.2 Implementation of nearest neighbor algorithm to increase data coverage

7

8 After implementing the prioritization rules there remain 5 arcmin cells with no carbon data
 9 coverage for a given land type and carbon pool. This is expected since the land cover data used
 10 to generate the carbon masks (ESA CCI land cover data) may be different from the land cover
 11 data used in HYDE, SAGE. We therefore implement a nearest neighbor algorithm to interpolate
 12 data to each ‘no data’ cell based on availability in 40 neighboring cells. This algorithm fills the
 13 target cell and land type with the corresponding carbon data of the closest cell with matching
 14 land type. If no matches are found within the prescribed window then the target cell remains
 15 without data for that particular carbon pool and that particular land type.

16

17 Carbon data coverage after interpolation is reasonable with the exception of a few land types.
 18 Table 4 shows the data coverage by land type after implementation of the nearest neighbor
 19 algorithm. All but three land types have over 80% data coverage for soil and vegetation carbon.



1 At least 25% of Tundra and Polar desert cells remain without carbon data. This is likely a result of
 2 differences in way Tundra land cover is defined by different datasets.
 3
 4 There have been more recent efforts to collect soil carbon data specifically for the permafrost and
 5 Tundra regions such as that by Hugelius et al.(Hugelius et al., 2014). This suggests that a future
 6 area of work would be to incorporate these more detailed datasets into either the source data or
 7 our processing workflow. Along with Tundra and Polar deserts, over 20% of the Urban land cells
 8 do not have carbon data. This is once again likely due to the different definitions of Urban land
 9 cover indifferent datasets. Our data coverage suggests that there exists more uncertainty in the
 10 Tundra, Polar, and Urban carbon values purely based on limited data availability. Recognizing
 11 and quantifying data availability by land type enables users to utilize their own judgement when
 12 using the carbon values for these land types.
 13

Land type	Total 5arcmin grid cells	Vegetation carbon Percentage unfound (NO DATA cells after interpolation)	Soil carbon Percentage unfound (NO DATA cells after interpolation)
Pasture	1195396	2.3	2.3
Cropland	952850	17.0	17.0
Grassland/Steppe	498404	15.0	14.6
OpenShrubland	274296	16.0	16.0
Desert	195579	1.0	1.1
TropicalEvergreenForest/Woodland	190780	0.0	0.3
Savanna	173776	8.0	7.6
BorealEvergreenForest/Woodland	148756	0.0	0.0
PolarDesert/rock/ice	132021	29.0	24.9
Urban	119597	22.3	22.3
TemperateDeciduousForest/Woodland	86922	1.0	1.1
DenseShrubland	78065	10.0	9.5
TemperateNeedleleafEvergreenForest/Woodland	71600	1.0	0.5
BorealDeciduousForest/Woodland	65824	0.0	0.4



TropicalDeciduousForest/Woodland	56377	1.0	1.4
Tundra	25000	29.0	24.9
TemperateBroadleafEvergreenForest/Woodland	14395	0.0	0.3

1

2 *Table 4: Details of NODATA cells after nearest neighbor interpolation*

3

4 **2.3 Stage 3 - Aggregating raster carbon data to 699 land regions**

5 As a final step, we pass the 72 rasters generated in Stage 2 and use those as inputs to the Moirai
 6 land data system. The land data system uses the inputs to aggregate the values to 699 land
 7 regions from the 5 arcmin grid cell level. The 699 land regions are the intersection of 235 water
 8 basins and 207 countries and are shown as a map in **SI figure 1**. The final carbon state values
 9 for each land type are aggregated to each land region for each carbon pool (aboveground
 10 biomass, belowground biomass, soil 0-30 cms). These outputs are available as a tabular text file.
 11 The moirai land data system performs this aggregation using the same land masks for the year
 12 2010 which are used in the Stage 2 processing. The basic aggregation performed by moirai is
 13 summarized in equation 3 below

14

$$15 \text{Carbon_tabular}_{pool, GLU, state, LT} =$$

$$16 \text{state} \left(\begin{array}{cc} \text{Carbon_5arcmin_LT}_{pool, j} & \text{Carbon_5arcmin_LT}_{pool, j+2} \\ \text{Carbon_5arcmin_LT}_{pool, j+1} & \text{Carbon_5arcmin_LT}_{pool, j+n} \end{array} \right) \quad (3)$$

17

18 Where,

19 pool is the carbon pool (aboveground biomass, belowground biomass, topsoil (0-30 cms)),
 20 state is the aggregation method (area-weighted average, median, min, max, q1, q3),

21 GLU represents a land region which is an intersection of 207 country boundaries and 235
 22 watershed boundaries,

23 j is the grid cell index for each 5 arcmin grid cell in a basin with land type LT,

24 n is the total number of cells in a basin for a given land type,

25 and LT is the land type.

26

27 **2.4 Stage 4 – Deriving any other percentile using our six statistical states**

28

29 Using our six summary states, users can calculate any percentile for the carbon value in any pixel
 30 for each of our 19 land types and three carbon pools (soil, above ground biomass, below ground
 31 biomass). These values can also be calculated directly for a land region/water basin. The
 32 percentile values can be calculated assuming that carbon values are lognormally distributed (this
 33 is established in our analysis below- See section 3.1) The steps to calculate any percentile are as
 34 follows,



- 1 1. Compute a mean value as a natural log of the median state. Since the distribution of carbon
2 values is lognormal, the natural log of our median would be an estimated mean for the
3 lognormal distribution.
- 4 2. Compute an estimated standard deviation using a natural log of the Q3 and the mean value
5 in step 1, specifically we use the formula- $(\text{LN}(Q3) - \text{LN}(\text{mean}))/0.675$.
- 6 3. Estimate the percentile value from the mean and standard deviation above. Since the logged
7 distribution is normal, users can compute this value using a z table for a normal distribution.
- 8 4. Calculate the exponent of the value in step 3.
- 9 5. Constrain this value to the max observed value in our dataset.

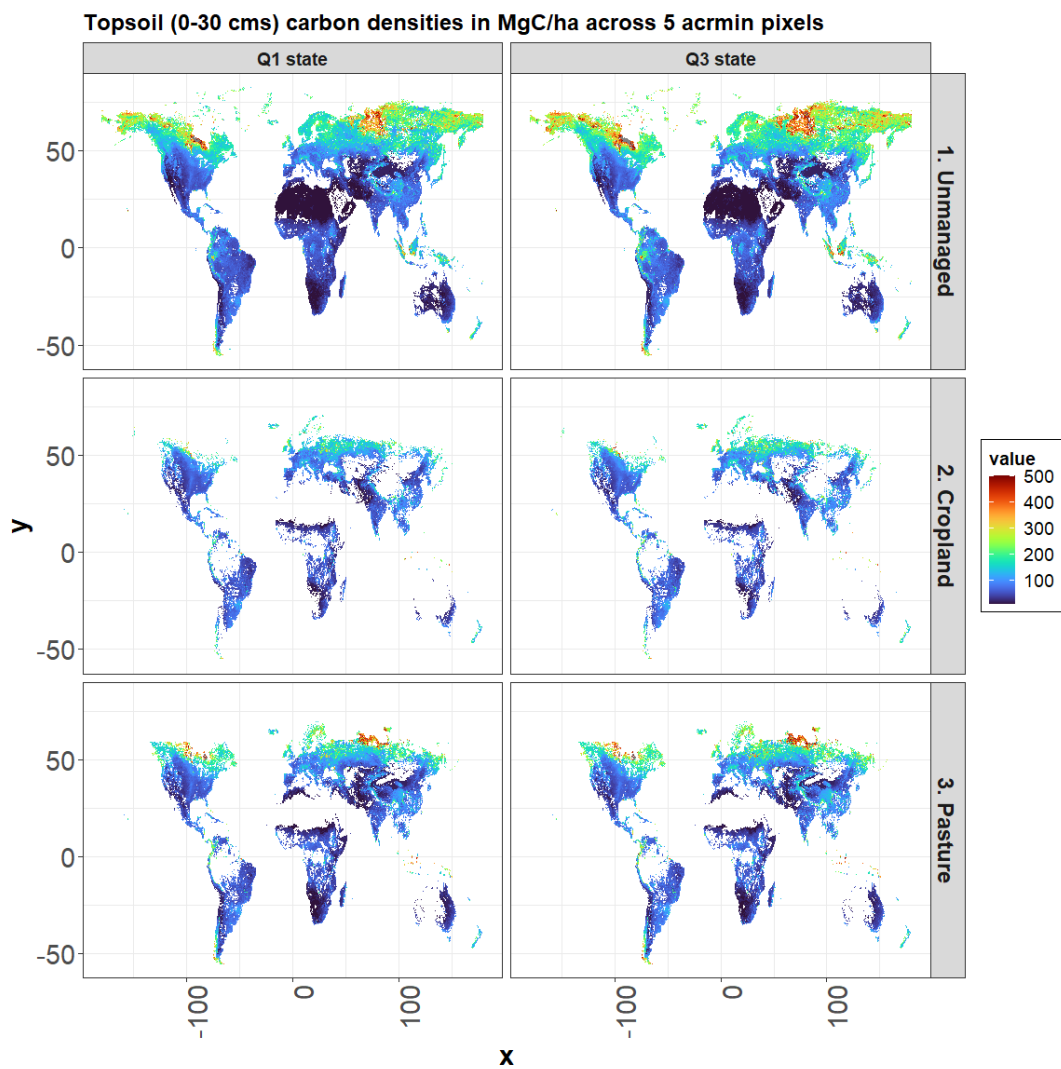
10 This method would enable a timely calculation of percentiles and would be much faster than re-
11 running the code to derive individual percentiles using re-sampling.

12

13 **3. Analysis, Uncertainty and data validation**

14

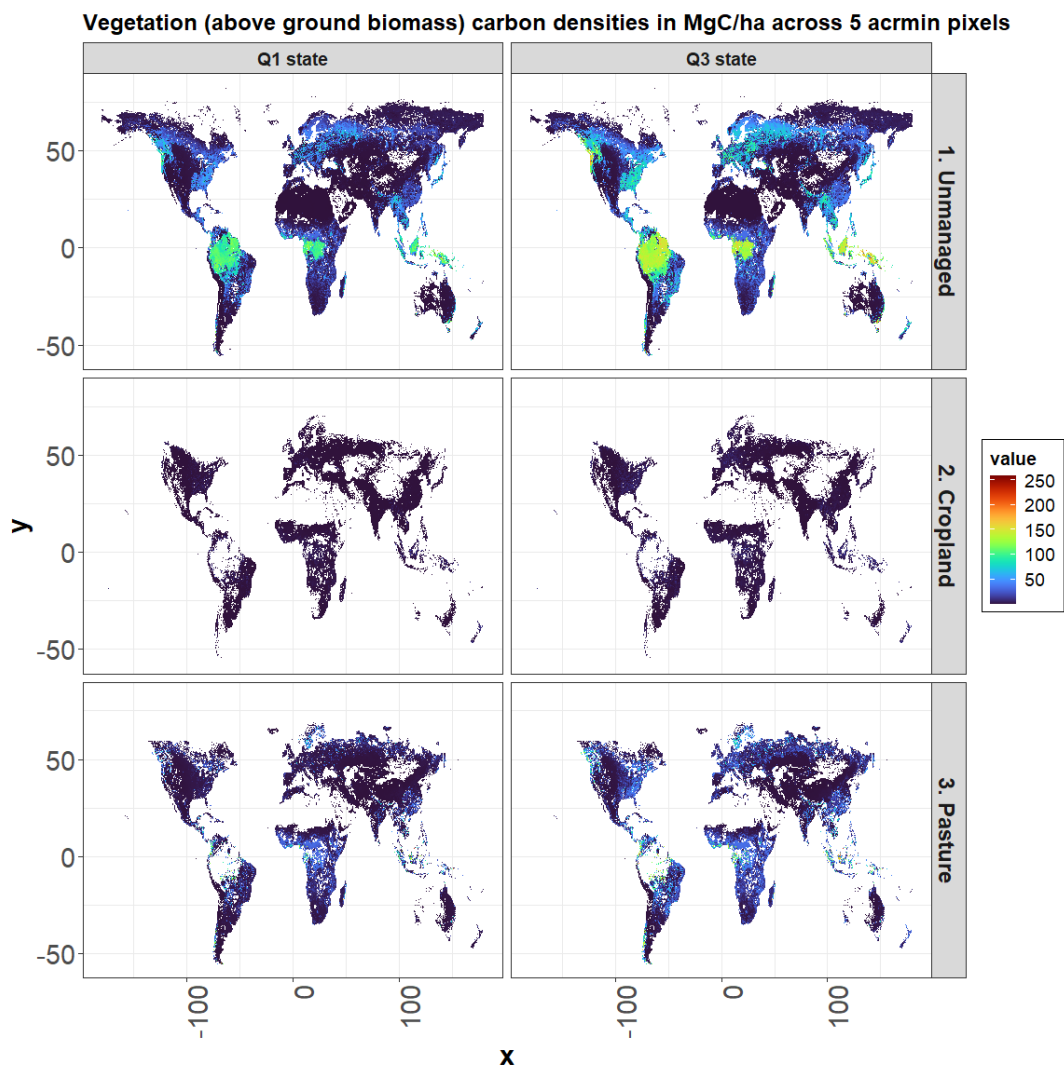
15 We first evaluate our main data products, namely the maps of soil and vegetation carbon
16 across gridcells by land type (e.g., Figure 2 and Figure 3), with the goal of identifying the
17 most appropriate carbon state for GCAM modeling, and then take a closer look at data
18 uncertainty and spatial variability. Note that the authors of the source data on soil (Hengl
19 et al.) and vegetation (Spawn et al.) did a detailed spatial validation of the data in their
20 respective papers. Our validation will focus on uncertainties that have been introduced
21 through our re-harmonization process. We will also compare our Q3 and 90th percentile
22 (determined as described above) estimates with similar estimates from the literature since
23 these estimates will be used to initialize GCAM.



1
2 *Figure 2: Soil carbon for each 5 arcmin grid cell in MgC/ha. Values are shown for two statistical states, namely the Q1 and the*
3 *Q3.*

4

5



1
 2 *Figure 3 Veg carbon (aboveground) in MgC/ha across 5 arcmin grid cells for aggregate land types for the Q3 state. Values are*
 3 *shown for two statistical states, namely the Q1 and the Q3*

4
 5
 6
 7
 8
 9
 10
 11
 12
 13

3A Comparison of harmonized carbon values to estimates of historical values

As mentioned above, this dataset was generated to initialize GCAM with spatially explicit carbon values for spin up and further simulation of the land system. This initialization requires the carbon values to represent a maximum potential carbon density because these values determine the limiting parameters for vegetation growth and soil carbon accumulation curves. The pre industrial carbon density has been estimated to be much higher than the contemporary carbon stored in land (Erb et al., 2018) due to a long



1 history of land use. Moreover, various studies have highlighted the difficulties with the
 2 calculation of the long term potential maximum(Fang et al., 2014). However, our
 3 development of the statistical states has allowed us to adopt a more systematic approach
 4 to selecting a data driven maximum value that we use to initialize GCAM.

5
 6 To select a data driven pre industrial equilibrium state, we compared the frequency
 7 distributions of carbon by pool within each land region for each land type with the final
 8 statistical states calculated. The frequency distributions represent a heterogeneous
 9 landscape at different stages of growth and management. The average or median values
 10 may be representative of the contemporary landscape, but not of an undisturbed
 11 landscape that has been allowed to equilibrate its carbon stocks. The maximum value in a
 12 land region may be an extreme outlier and likewise would not be representative of the
 13 undisturbed landscape. Our goal then is to find a value in between the contemporary
 14 average and the maximum that is representative of a long term potential maximum value.
 15 Fortunately, most distributions of soil carbon generally follow a log-normal shape with a
 16 long tail. For example, we analyzed the distribution of soil carbon in the Amazon basin
 17 (Figure 4) for different land types.

18
 19 One possible option for initialization is the Q3 statistical state. The Q3 statistical state
 20 value from these distributions does fall between the average and the maximum, as
 21 expected. Given the lognormal shape, the observations above the Q3 value are infrequent
 22 and can stretch to extremely high values. We also find that most vegetation carbon
 23 distributions follow a log-normal shape within each basin for each land type. However,
 24 forests have distributions that are more bimodal (Figure 5). Nonetheless, in these
 25 distributions the Q3 value provides an estimate of carbon that is reasonably higher than
 26 the contemporary average or median value. Using the Q3 values to initialize GCAM, we
 27 found that the initial carbon stock in the year 1700 would be approx 2144 PgC (1553 PgC
 28 of carbon in top soil and 591 PgC of vegetation biomass). This estimate is still on the
 29 lower end of other similar estimates from Walker et al. ,Erb et al, Houghton and
 30 Sanderman et al. Table 5 below summarizes our initial terrestrial carbon stock in 1700
 31 calculated from different sources-

Data source	Topsoil (0-30 cms) carbon in PgC	vegetation (above+ below ground biomass) in PgC
Erb et al 2019		916
moirai (Q3)	1553	591
moirai (90th percentile)	2063	966
Walker 2022		795
Sanderman et al. 2017	2119	43
Houghton 1999	1462	662



1 *Table 5: Initial potential terrestrial carbon stock calculated from different sources.*
2 *Sources from moirai are calculated using land maps in 1700. Sanderman et al.*
3 *represents a carbon stock in 1800 given no land use. Walker and Erb et al. are based on*
4 *potential vegetation carbon estimations.*

5
6 In addition to the Q3 value, we also use the estimated 90th percentile state in order to
7 represent a higher initialization of carbon in 1700. This 90th percentile is estimated from
8 our six summary states using the methodology outlined in section 2.4. This 90th percentile
9 provides an initial carbon stock of 3028 PgC (2063 PgC of carbon from topsoil and 966
10 PgC of vegetation biomass). Using these two states for initialization helps us understand
11 the sensitivity of the model to the initial value.

12
13 In the case of forests, we note that we derive carbon values for Forests as a whole and do
14 not differentiate between Primary Forests and Secondary Forests. This is a result of lack
15 of availability of fine resolution land masks that differentiate between primary and
16 secondary forests. This likely means that our long term potential maximum forest carbon
17 densities include the impact of harvesting especially in regions with high levels of forest
18 harvests. As more fine resolution data on different types of forests become available, a
19 logical next step would be to derive separate carbon densities for this particular land type.

20

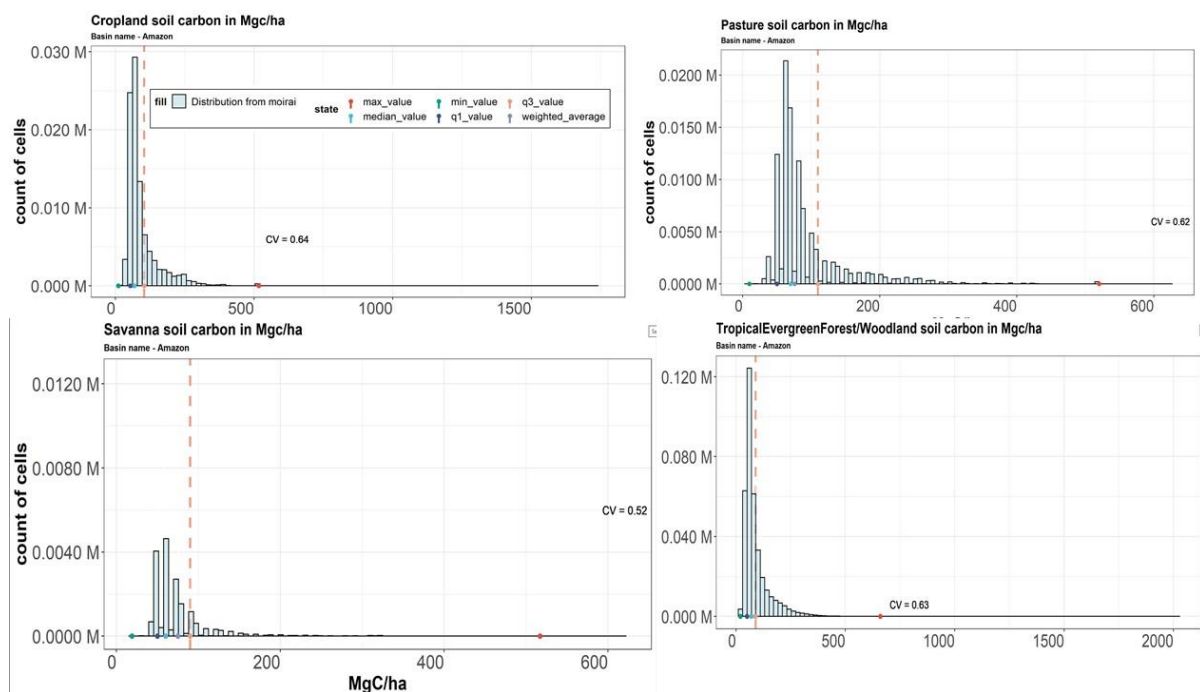
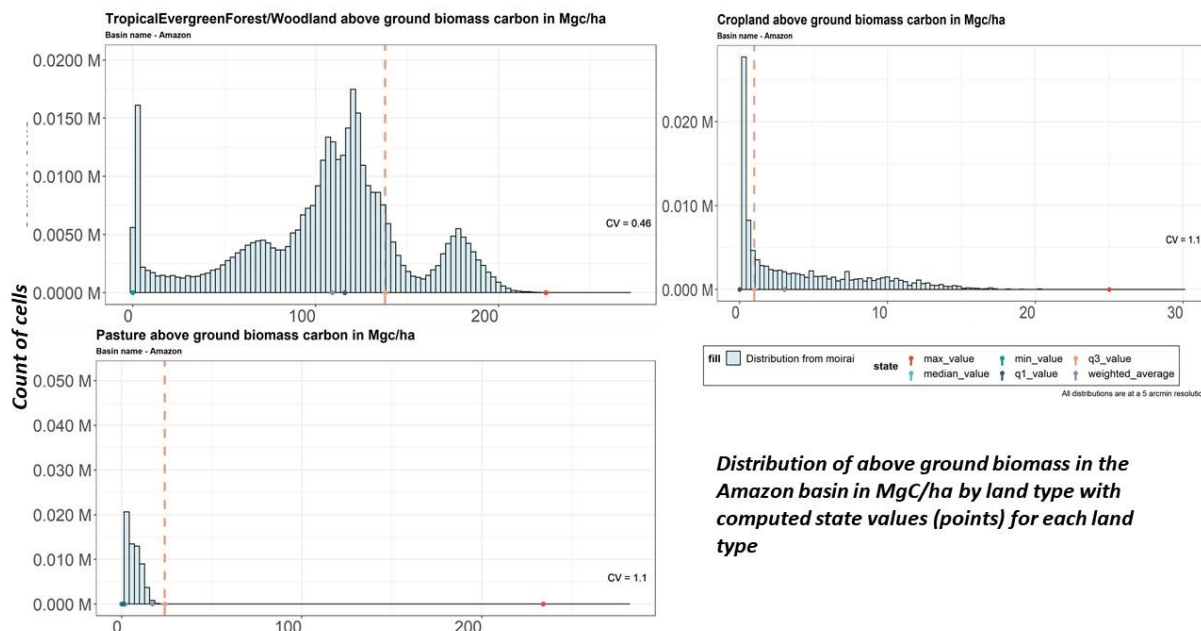


Figure 4: Within basin distributions of soil carbon in MgC/ha for the Amazon basin. Each facet shows a distribution for a land type. The final basin level statistical states are shown as dots with the Q3 state shown as the orange line.



Distribution of above ground biomass in the Amazon basin in MgC/ha by land type with computed state values (points) for each land type

Figure 5: Within basin distributions of aboveground biomass carbon in MgC/ha for the Amazon basin. Each facet shows a distribution for a land type. The final basin level statistical states are shown as dots with the Q3 state shown as the orange line.



1 **3B Comparison of carbon values to other estimates of long term potential carbon by**
2 **grid cell**

3
4 Since the Q3 and 90th percentile values in our dataset will be used to represent pre-
5 industrial carbon densities in GCAM, we compared our values to similar estimates in the
6 literature. Specifically we compare the 90th percentile values at the pixel level with other
7 estimates, since this statistical state represents a high carbon initialization compared to the
8 Q3 and produces a global terrestrial carbon stock that is in line with other estimates of
9 potential terrestrial carbon.

10
11 Sanderman et al. 2017 generated a pre-industrial soil carbon map for top soil in the year
12 1800. This map assumed no land use in that year. Similarly Walker et al. (Walker et al.,
13 2022) generated a similar map for potential carbon in above and below ground vegetation.
14 For a valid comparison we compared only our unmanaged land carbon values with these
15 estimates (Figure 6 and 7).

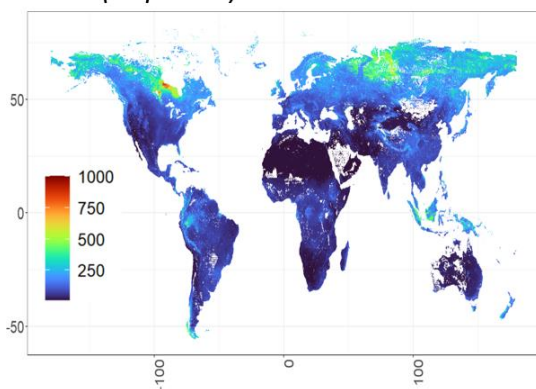
16
17 We found that in the case of soil carbon, even though our maps track well with the maps
18 from Sanderman et al. in terms of the overall spatial distribution, the mean error (moirai
19 90th percentile – Sanderman et al.) across gridcells that is close to -23%. There are some
20 upper latitude pixels from the Sanderman et al. dataset that show almost 100% higher
21 values compared to our data.

22
23 In case of aboveground vegetation carbon, the mean percen error (moirai 90th percentile –
24 Walker et al. 2022) is -17%, which is lower than for soil carbon. The largest errors were
25 observed for forest pixels. This is likely due to the combination of Primary and Secondary
26 forests into a single forest category in our dataset (as described above). The highest
27 differences between datasets are observed in forest pixels with high level of forest
28 harvesting (Central and West Africa and South and East Asia).

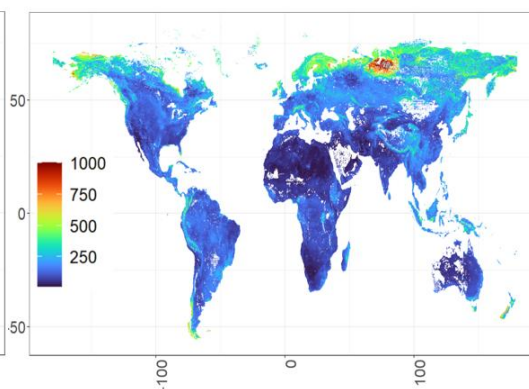
29
30



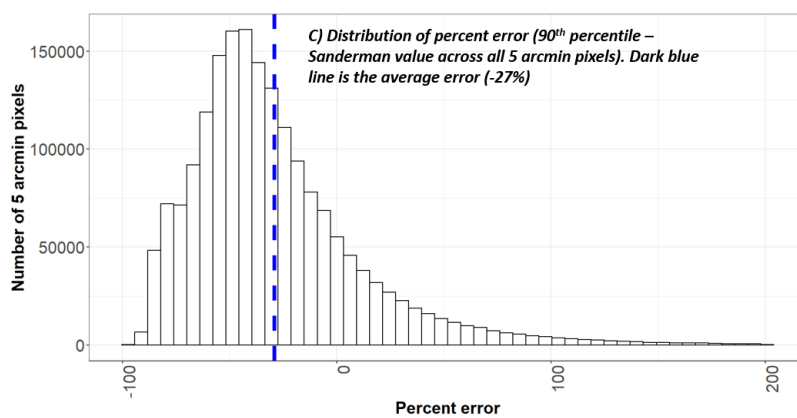
A) Soil carbon (0-30 cms) in MgC/ha at 5 arcmin for moirai (90th percentile)



B) Potential Soil carbon (0-30 cms) in MgC/ha at 5 arcmin for Sanderman et al. 2017



1
2

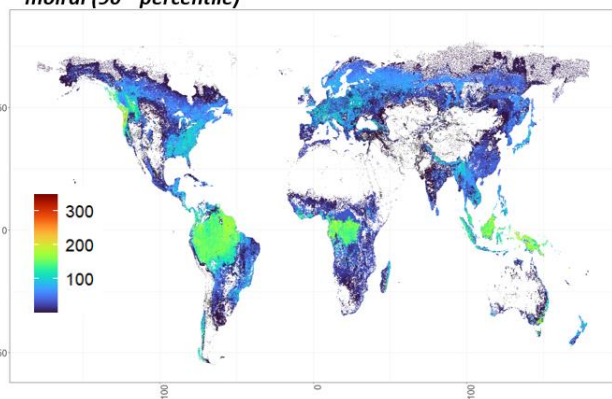


3
4
5
6
7

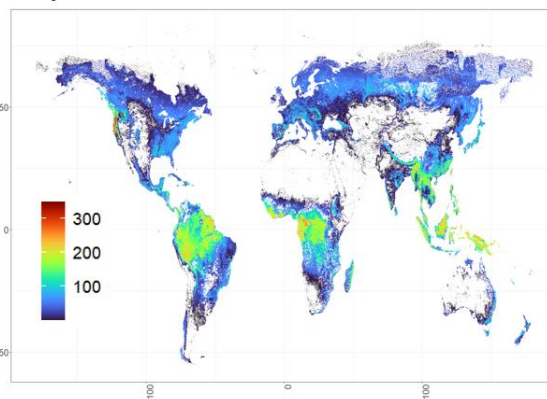
Figure 6: A.) Topsoil (0-30 cms) carbon in MgC/ha for 5 arcmin pixels using moirai 90th percentile B.) Top soil (0-30cms) carbon in MgC/ha from Sanderman et al. assuming a no land use condition. C.) Histogram showing percent error between A and B. Dark blue dashed line represents mean error across all pixels which is at -27%.



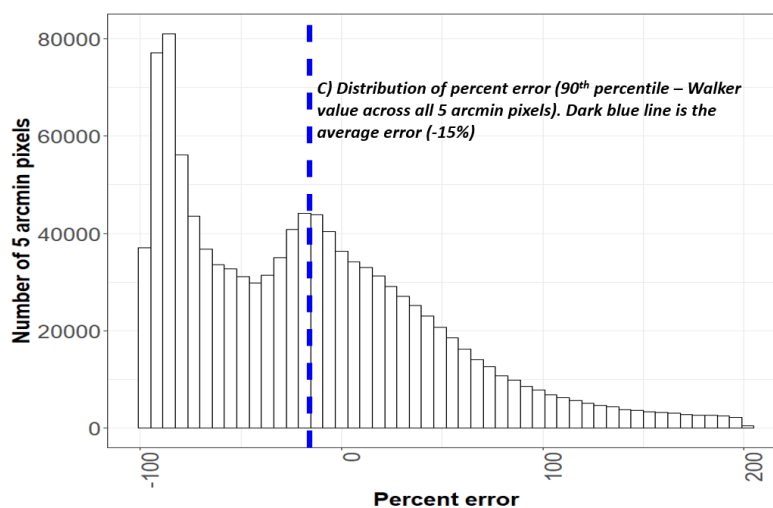
A) Above ground biomass in MgC/ha at 5 arcmin for moirai (90th percentile)



B) Potential aboveground vegetation in MgC/ha at 5 arcmin for Walker et al. 2022



1



2

3

4

5

Figure 7: A.) Vegetation (aboveground) carbon in MgC/ha for 5 arcmin pixels using moirai 90th percentile B.) Vegetation (aboveground) carbon in MgC/ha from Walker et al. constrained for initial land use C.) Histogram showing percent error between A and B. Dark blue dashed line represents mean error across all pixels which is at -15%.

6

7

8

9

10

11

12

13

14

15

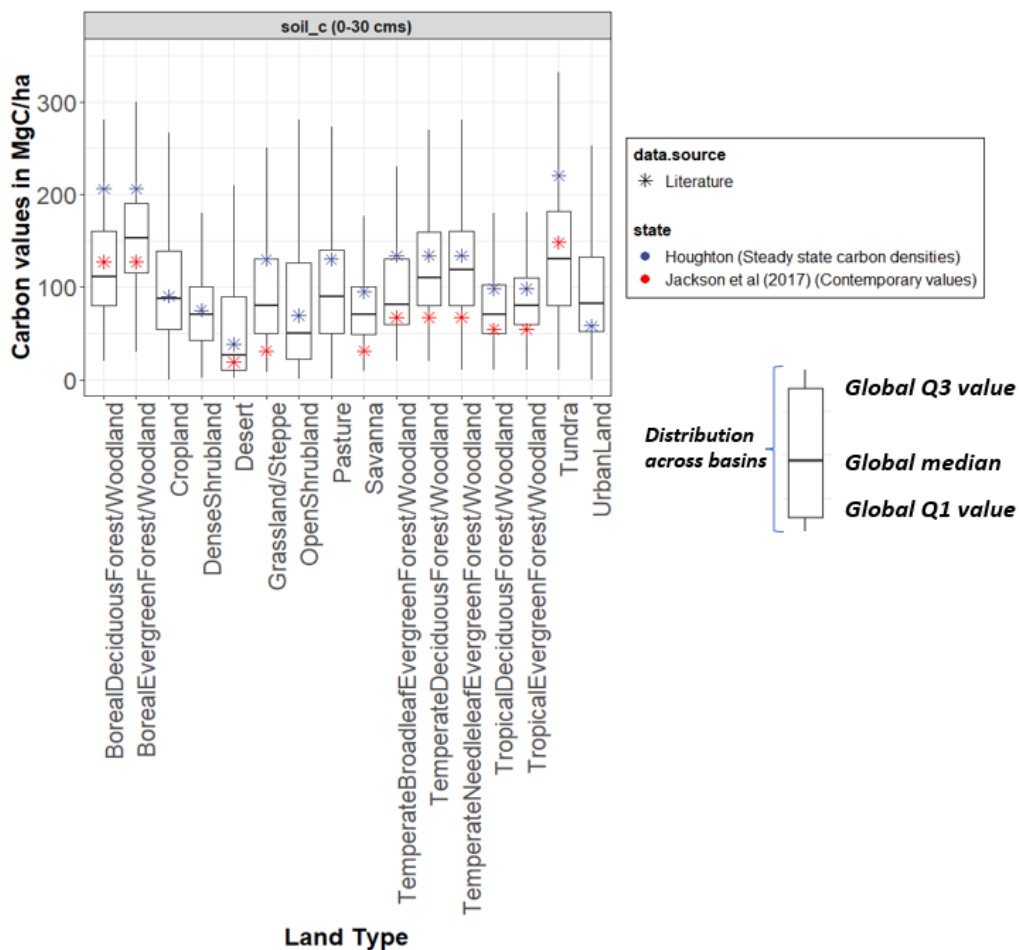
16

3C Comparison to C values to previously used in GCAM by land type and aggregate contemporary estimates

We compared the moirai densities by land type globally with similar carbon densities from Houghton (1999) (See **SI Table 2** and Figure 8 for a comparison of soil carbon estimates and **SI Table 3** and Figure 9 for aboveground biomass comparison). The Houghton carbon densities represent carbon values on undisturbed land differentiated by biome. We specifically compare against the Houghton carbon densities since those values were previously used for the spin up in GCAM. We also compared our statistical states with contemporary values where available (e.g. Jackson et al. for soil carbon and Vlek et al. for contemporary vegetation carbon).

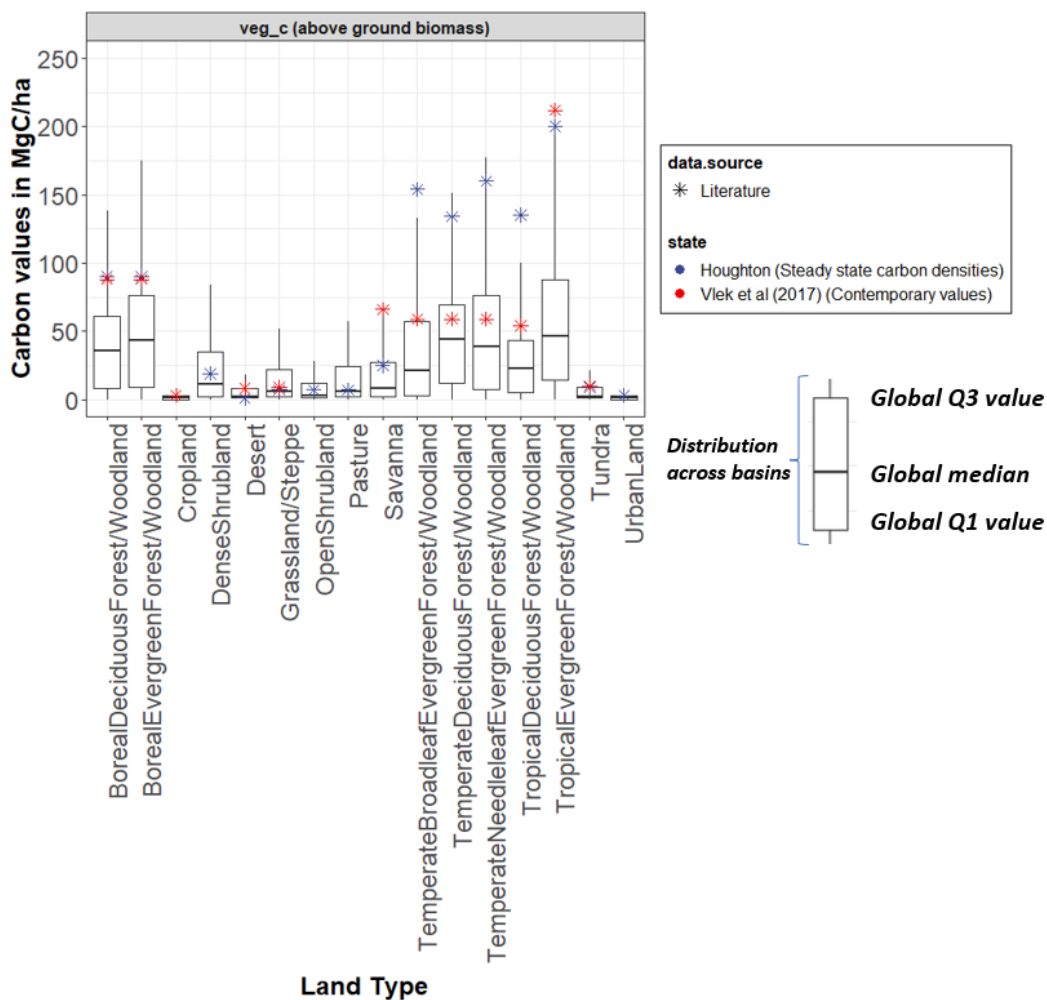


1
2 For soil carbon (Figure 8), we found that our values (Q3, 90th percentile) are generally
3 higher than the Houghton values globally for most land types. The values are especially
4 higher for Shrublands which are located in Boreal regions where the difference is approx
5 80 Mgc/ha. This is likely because the SoilGrids dataset shows high carbon values at high
6 latitudes and includes peat soils in its estimates (e.g., figure 3). The high values of soil
7 carbon at the upper latitudes may also be driven by high levels of predicted bulk density
8 at those locations(Tifafi et al., 2018).Another version of soil grids has recently produced
9 lower values in these regions to reduce the effects of peat soil estimates(Poggio et al.,
10 2021).
11
12 For cropland, our Q3 estimates of carbon are as high as forest soil carbon. This is
13 investigated in more detail in the sections below. Similarly the soil carbon under Urban
14 land cover is extremely high. This is likely due to how the samples were collected for
15 Urban land cover (these samples are collected in parks as opposed to built up areas). As
16 expected, the values in our range are higher than the contemporary values from Jackson
17 et al., especially in the Boreal regions. However, the Q1 values from our range are closest
18 to contemporary values for soils.
19
20 For vegetation carbon (Figure 9), in the case of forests, the carbon densities are
21 significantly scaled down across moirai states when compared to the literature (Houghton
22 for the pre-industrial values and Vlek et al. for contemporary values). This is not
23 surprising since in the Houghton inventory numbers, the spatial distribution of forest
24 carbon is an unknown especially for tropical forests (Houghton, 2005). Also, as noted
25 above our moirai values for forest carbon densities are a combination of Primary and
26 Secondary Forests and are therefore underestimate the long term potential maximum
27 carbon that can be stored in forests.
28
29 For grasslands and pastures, the moirai Q3 and 90th percentile estimate is higher than the
30 literature values. However, for this land type, the overall distribution of carbon values is
31 not very dispersed across basins. For forests, there is significant variation in values across
32 basins , likely since some forests may be more intensively managed or harvested
33 compared to others. Note that the reduction in carbon values compared to the global
34 values from Houghton will likely reduce the afforestation response in GCAM under
35 scenarios where carbon in forests is priced.
36
37
38



1
 2 Figure 8: Carbon densities globally by land type across basins and Houghton pre-industrial state. Also shown is the
 3 contemporary carbon density value by land type from Jackson et al.

4
 5



1
 2
 3
 4
 5
 6
 7
 8
 9
 10
 11
 12
 13
 14
 15

Figure 9: Carbon densities globally by land type for vegetation carbon across basins and Houghton values (Undisturbed carbon in 1850). Also shown is the contemporary carbon density value by land type from Vlek et al.



1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45

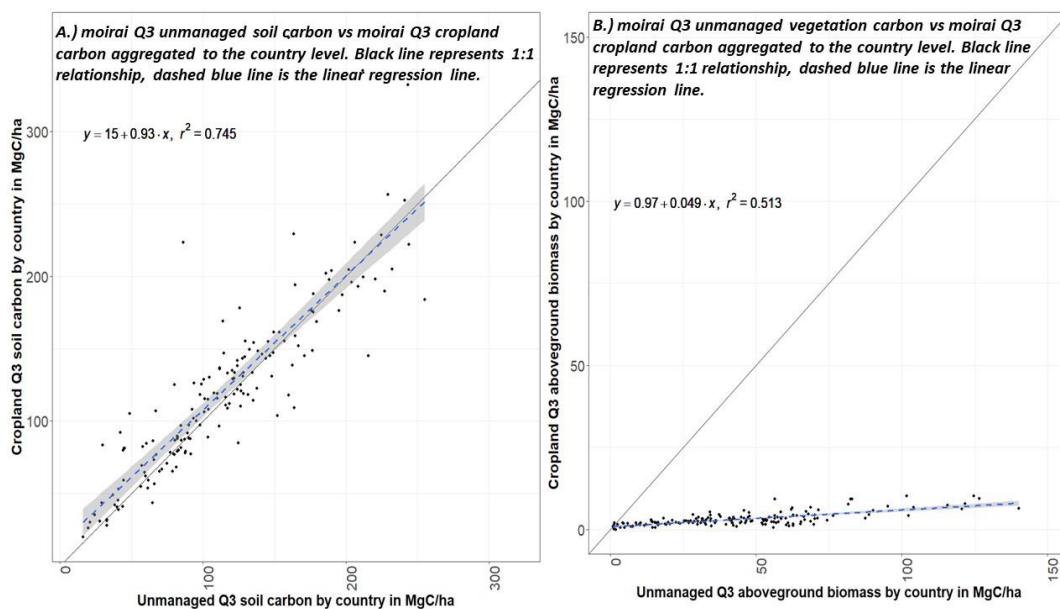
3.1 Uncertainties in re-harmonized carbon data (spatially and across land types)

Here we explore uncertainty in the available data by further examining spatial distributions, aggregation statistics, and land type considerations.

i.) Do managed land types show a depreciation in carbon compared to unmanaged land types?

Studies show that managed land (i.e., Cropland, Pasture, Urbanland) has depleted carbon stocks in relation to undisturbed land (Cooper et al., 2021; Sanderman et al., 2017; Wei et al., 2014). The aim of processing the spatial managed land carbon data and adding it to Moirai is to obtain contemporary estimates for these lands that can be used in modeling rather than assuming a global value or that managed lands have a fixed fraction of unmanaged land carbon. Carbon data values do not correspond with a long term potential maximum for these managed land types by definition, as these land are actively disturbed. However, we still want higher than average carbon values for the parameters that define the limits of carbon accumulation for these land types. We expect that the carbon data reflect the effects of these managed land types and that our desired values would be lower than those for the surrounding unmanaged land types. We checked this expectation by first comparing Q3 carbon values for soil and aboveground biomass for Cropland with the corresponding values for Unmanaged land cover in each of our land regions (Figure 10). We found that Cropland soil carbon values do not show a consistent depletion for soil carbon compared to Unmanaged land. The reason for these differences among carbon pools is rooted in the source data sampling and processing methodologies. In case of the soilgrids dataset, the authors state that cropland soil carbon samples were largely collected in the US. In case of the vegetation carbon dataset from Spawn et al., the vegetation carbon was calculated for each crop type based on yields which explains the low values on cropland compared to unmanaged land.

For cropland, yields are determined from harvested area and production data, while the carbon data are used for land use emissions and when valuing carbon in forcing scenarios. To address the relatively high cropland soil carbon data in our modeling experiments we reduce these data by 30% before using them in GCAM. Previous studies have found a similar loss of soil carbon through agricultural practices and land conversion from unmanaged land types to cropland (Cooper et al., 2021; Wei et al., 2014).



1
 2 *Figure 10: Comparison of unmanaged Q3 carbon densities and Cropland carbon densities for A.) soil carbon and B.)*
 3 *aboveground biomass. Values here are aggregated to individual countries*

4
 5 We performed a similar analysis for pasture carbon densities and found that
 6 Pasture carbon shows depletion or lower values compared to unmanaged land
 7 cover both for both soil and vegetation. In this case, it is reasonable to use the Q3
 8 soil and vegetation carbon values for Pasture in GCAM without adjustment.

9
 10
 11 ii.) *Assesing spatial uncertainties in soil and vegetation carbon within and across*
 12 *basins*

13
 14 We have established that carbon distributions within land regions generally follow a
 15 lognormal pattern for soil carbon and for vegetation carbon for most land types
 16 while vegetation carbon for Forests has a more bimodal distribution. However,
 17 there may be more dispersion across values in some basins for some land types
 18 compared to others. To assess this systematically, we computed a quartile
 19 coefficient of dispersion (QCD) for each basin and land type as:

20
 21
$$QCD_{GLU,LT,pool} = (Q3_{GLU,LT,pool} - Q1_{GLU,LT,pool}) / (Q3_{GLU,LT,pool} +$$

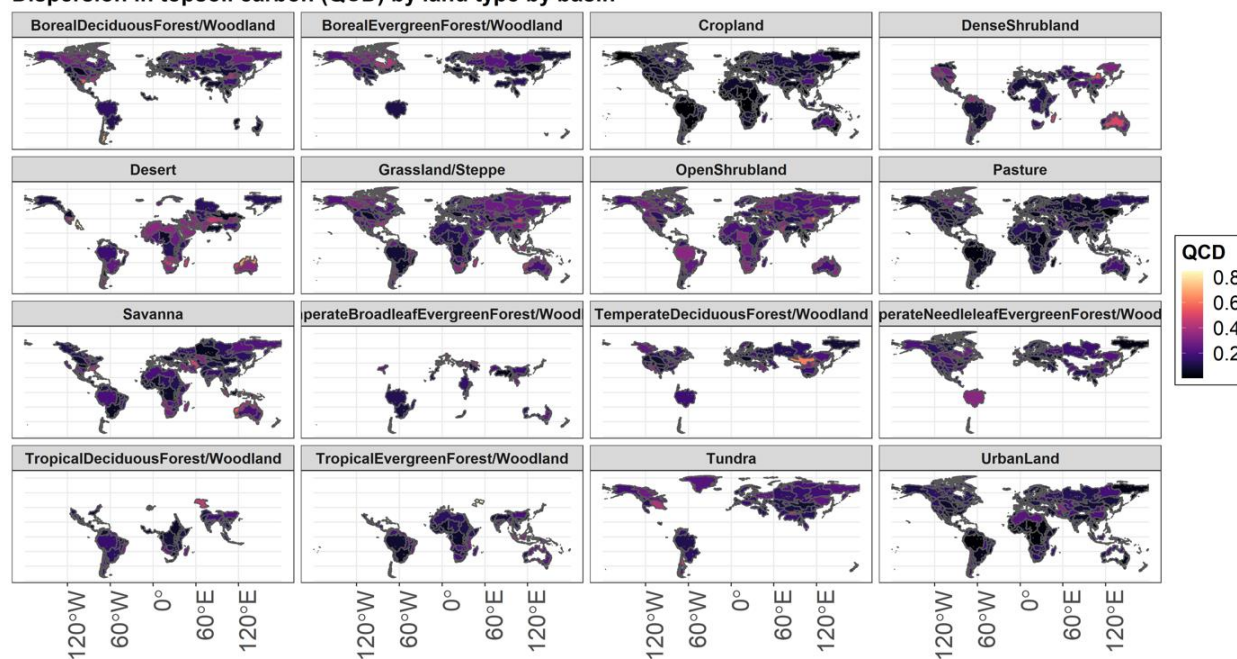
 22
$$Q1_{GLU,LT,pool}) \quad (5)$$

23
 24 Where,
 25 pool is the carbon pool (aboveground biomass, belowground biomass, topsoil (0-30 cms)),
 26 GLU represents a land region which is an intersection of basin boundaries and country
 27 boundaries,
 28 and LT is the land type.



1
2
3 The QCD values ranges from 0 to 1 where a value towards zero indicates less dispersion within a
4 region-land type-carbon pool combination and a value towards 1 indicates more dispersion.
5
6 The QCD values for soil carbon (Figure 11) are generally similar across most basins across land
7 types. This is expected since the distributions of soil carbon are generally lognormal. However, in
8 some basins the QCD value is consistently high and similar across land types. This mainly occurs
9 in individual basins in Russia and Indonesia which have high levels of peat soils which would
10 mean that the level of dispersion across cells would be high since some cells would contain peat
11 soils whereas others would not.
12

Dispersion in topsoil carbon (QCD) by land type by basin



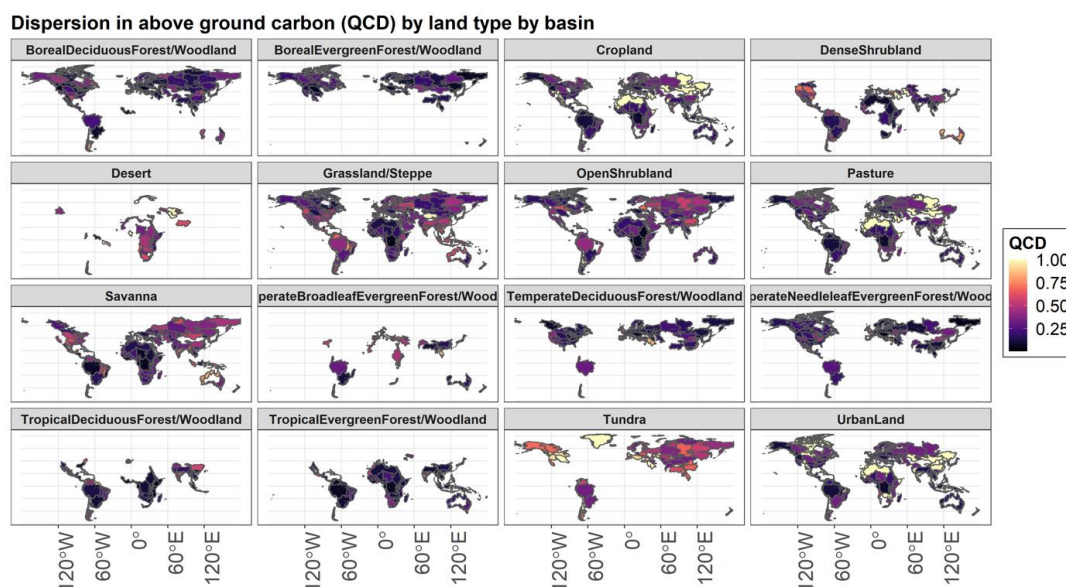
13
14
15
16
17
18
19
20
21
22
Figure 11: QCD values for topsoil carbon across basins and land types

14
15 Based on QCD values across basins and land types for vegetation carbon (Figure 12), we
16 observe that there is significant variation in the QCD values within and across basins for Tundra
17 (with values ranging from 0-1). This is likely due to the way Tundra pixels are defined in our
18 dataset (they encompass different vegetation types). Similarly, there is significant variation
19 within and across basins for grasslands, savannah and pastures, which is once again likely due to
20 the definitions of what constitutes grasslands in the base land cover dataset. While there are also
21 variations in vegetation carbon values for cropland and urbanland, the overall range of values for
22 these land types when it comes to vegetation carbon is low (Figure 7). QCD values for forests



1 across and between basins is lower. This may be due to the more narrow definitions for what
 2 constitutes forests across datasets.

3
 4



5
 6 *Figure 12: QCD values for aboveground vegetation across basins and land type*

7
 8
 9

4. Results from implementation of spatially explicit carbon in GCAM

10 We make the following assumptions when implementing the carbon densities in GCAM,
 11 based on the analyses above:

- 12 a.) The Q3 carbon values and the 90th percentile values are used throughout to reflect two
- 13 potential options for a long term potential maximum state of carbon in 1700
- 14 b.) Cropland soil carbon is reduced by a factor of 0.3 (30% reduction) for all basins to
- 15 reflect the effects of management
- 16 c.) Tundra, Urban, Desert, and Polar desert/rock/ice do not change in GCAM and so the
- 17 assigned carbon values do not influence model simulations. If a model does include
- 18 dynamics for these land types, then the associated uncertainties should be addressed.

19
 20
 21

4.1 Results from historical spin up

22
 23
 24
 25
 26
 27
 28

We initialize GCAM using our two options identified above. This results in a pre-spin up carbon stock of 2144 PgC (1553 PgC in soil and 591.7 PgC in vegetation) when using the Q3 state and a carbon stock of 3028 PgC (2063 PgC in top soil and 965 PgC in vegetation) when using the 90th percentile. Note that these initialization values are calculated using the land cover in 1700, which does include some managed area, and the spatially explicit carbon.



1 During the spin up (Table 6), this carbon is reduced to 1967 PgC in 2015 when using the
2 Q3 state(1481 PgC of topsoil carbon and 486 PgC of vegetation carbon) as a result of
3 historical land transitions. Similarly during the spin up, this carbon is reduced to 2758 PgC
4 when using the 90th percentile values (1965 PgC of topsoil carbon and 793 PgC of
5 vegetation carbon).
6

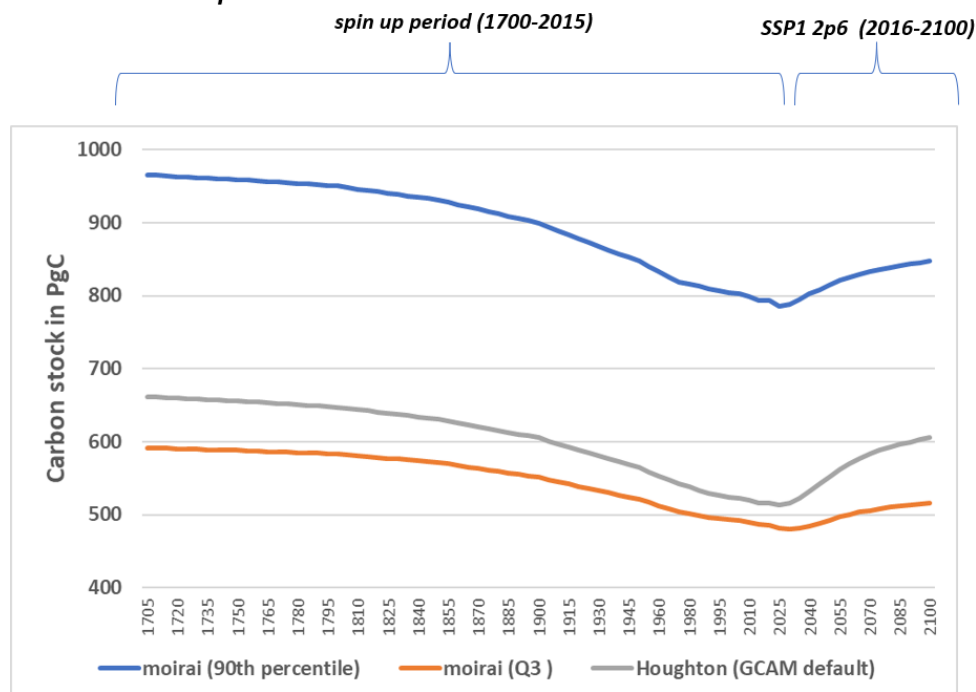
7 An important point to note is that while the 90th percentile generates results more in line
8 with independent pre-industrial estimates (e.g. Walker Sanderman, Erb), the Q3 state
9 results in more realistic contemporary values in 2015 during the GCAM spin up. For
10 example, the Q3 state results in a contemporary value of 486 PgC of vegetation carbon in
11 2015, which is closer to contemporary vegetation carbon stock estimates. Whereaes, the
12 90th percentile results in a global vegetation carbon stock of 793 PgC. Using the 90th
13 percentile would effectively result in an unrealistically high initial vegetation carbon stock
14 that is close to equilibrium in 2015. Furthermore, when running a carbon price scenario
15 using the 90th percentile densities (Figure 13), the model would add another 54 PgC of
16 vegetation carbon through afforestation that would result in a very unrealistic value of
17 vegetation carbon in 2100- Close to 847 PgC which is higher than undistrurbed carbon
18 stocks in 1700. When using the Q3 densities, this vegetation carbon stock in 2100 is close
19 to 515 PgC (Additional 30 PgC of carbon added through planted forests).
20

21 Another point to note is that the amount of global historical emissions (1700-2015)
22 produced by the Q3 initialization is 176 PgC which is much lower than the global
23 historical emissions using the 90th percentile of 270 PgC. For context, the Global Carbon
24 Project (as of 2021) produced an estimate of annual LUC emissions from 1700-2015 of
25 196 PgC(Friedlingstein et al., 2022). Figure 14 below shows the LUC emissions for the
26 historical periods for the GCP and our two initialization options. As seen in the figure, the
27 90th percentile produces consistently higher annual LUC emissions.
28

29 Given the above results from the spin up, we found that the Q3 value from our dataset is
30 appropriate for initialization and use in GCAM when using the model to estimate
31 contemporary C dynamics. While the 90th percentile better resembles independent
32 estimates of pre-settlmenet stocks, it results in substantial overestimation when used to
33 estimate contemporary C fluxes. This is a result of assumptions and processes within
34 GCAM pertaining to carbon dynamics. As such, What is appropriate for other models
35 would likely be different and would require a similar analysis.
36
37



Example of spin up results in GCAM for vegetation carbon using several initialization options



1
2

Figure 13: Descriptions of the results of the spin up process. Global vegetation carbon during spin up (1700-2015) and the SSP1 2p6 climate forcing scenario (2016-2100) for our initialization options.

3
4

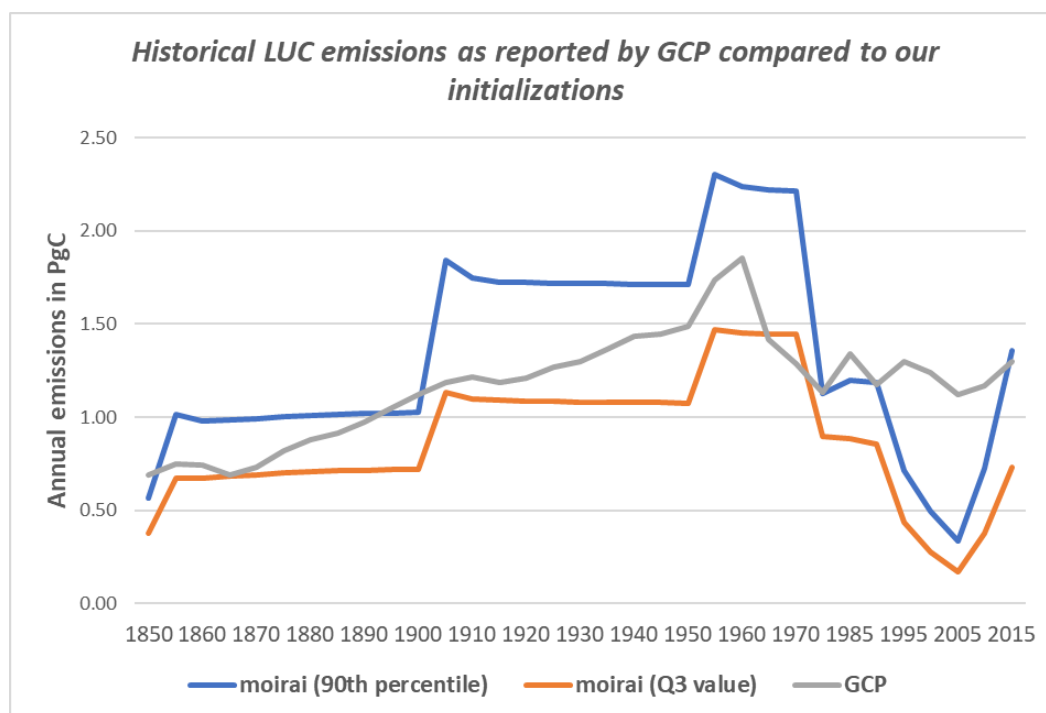
Table 6: Results from the historical spin up

Initialization	carbon pool	Initial value in PgC In the year 1700)	Contemporary value after spin up in PgC (2015)	Historical emissions (PgC) (between 1700 and 2015)	Value in 2100 under SSP1 2p6	Additional carbon sequestered during afforestation scenario (2100 value- 2015 value)
Houghton	vegetation carbon	662.0	516.1	145.9	605.3	89.2
moirai (Q3 value)	vegetation carbon	591.7	486.3	105.4	515.9	29.6



moirai (90th percentile)	vegetation carbon	965.8	793.2	172.6	847.9	54.7
Houghton	soil carbon (top-soil)	1243.5	1181.4	62.1	1220.0	38.6
moirai (Q3 value)	soil carbon (top-soil)	1320.5	1249.1	71.4	1274.6	25.5
moirai (90th percentile)	soil carbon (top-soil)	1753.0	1655.2	97.8	1700.6	45.4
Houghton	Total terrestrial carbon	1905.5	1697.5	208.0	1825.3	127.8
moirai (Q3 value)	Total terrestrial carbon	1912.3	1735.4	176.9	1790.6	55.2
moirai (90th percentile)	Total terrestrial carbon	2718.8	2448.3	270.5	2548.4	100.1

1
 2



1
2 Figure 14: Annual Global LUC emissions from GCP 2021 and our two initialization options

3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24

4.2 Results from climate forcing scenario

We use one climate forcing socioeconomic scenario (SSP1 2p6) with a maximum radiative forcing level of 2.6 watts per square meter by 2100 to assess how the new carbon data influence land projection in GCAM. Under this scenario land carbon prices are implemented to assign value to terrestrial carbon at the same rate as carbon is valued in the energy system. GCAM by default uses carbon densities from Houghton et al. (1999) which are described in SI Table 5 (soil) and SI Table 6 (vegetation). Note that the changes in land cover under the climate forcing scenario are driven by relative levels of carbon across land types rather than absolute levels of carbon. Therefore, even if forest carbon in some tropical regions are lower than other estimates, forests still sequester much more carbon compared to other land types in these regions. Below, we will compare results for the climate forcing scenario when using values from Houghton, the moirai Q3 value and the moirai 90th percentile.

The global land allocation comparison under SSP1 2p6 scenario in GCAM (Figure 15) shows that the afforestation/reforestation response is greatly reduced as a result of the spatially explicit carbon (the increase in forest cover from 2020 to 2100 globally is only 3.2 million km² when using the moirai Q3 as opposed to 7 million km² with the



1 Houghton carbon). IAMs (Including GCAM) generally show a very optimistic
2 afforestation response for this scenario that ranges from 0.5 to 12 million km² of trees
3 planted as part of a nature based carbon sequestration strategy under SSP1 2p6 (Popp et
4 al., 2017). The afforestation response in IAMs has been considered too optimistic in some
5 studies (Pongratz et al., 2021). The reduced forest expansion in GCAM with the new
6 carbon data is largely driven by lower vegetation carbon densities in the new data. When
7 using the 90th percentile, the afforestation values are the lowest at 0.1 million km². This
8 is expected given that the 90th percentile carbon values are much higher so much lower
9 increases in forest cover are required to meet additional afforestation targets. Usage of the
10 90th percentile values leads to minimal changes across land types under a climate forcing
11 scenario. This would be expected since the carbon pool generated by this scenario in
12 2015 is already close to equilibrium. This further illustrates why the Q3 state from our
13 dataset is a better choice for initialization in a model like GCAM.

14
15 However, the afforestation responses are diverse by regions. In case of tropical forests,
16 there is an increase in the afforestation response with our updated carbon densities. In
17 case of boreal forests, the opposite is true.

18
19 Global Cropland and Shrubland dynamics show a more complicated response. The
20 reduced emphasis on forest expansion reduces the need for Cropland abandonment.
21 Cropland also sequesters more soil carbon in some regions (even with the 30% reduction
22 factor), which also reduces abandonment. This Shrubland response is also enhanced by
23 higher Shrubland vegetation carbon densities.

24
25 Regional responses are dictated by their respective land type distributions (e.g., SI Figure
26 5). For example, in Russia the afforestation strategy is completely replaced with a
27 shrubland and grassland preservation strategy. This is expected since the region has a
28 relatively high amount of boreal forests. In South Asia however, where non-forest land
29 types dominate, forest expansion persists and is supplemented by shrubland expansion.
30

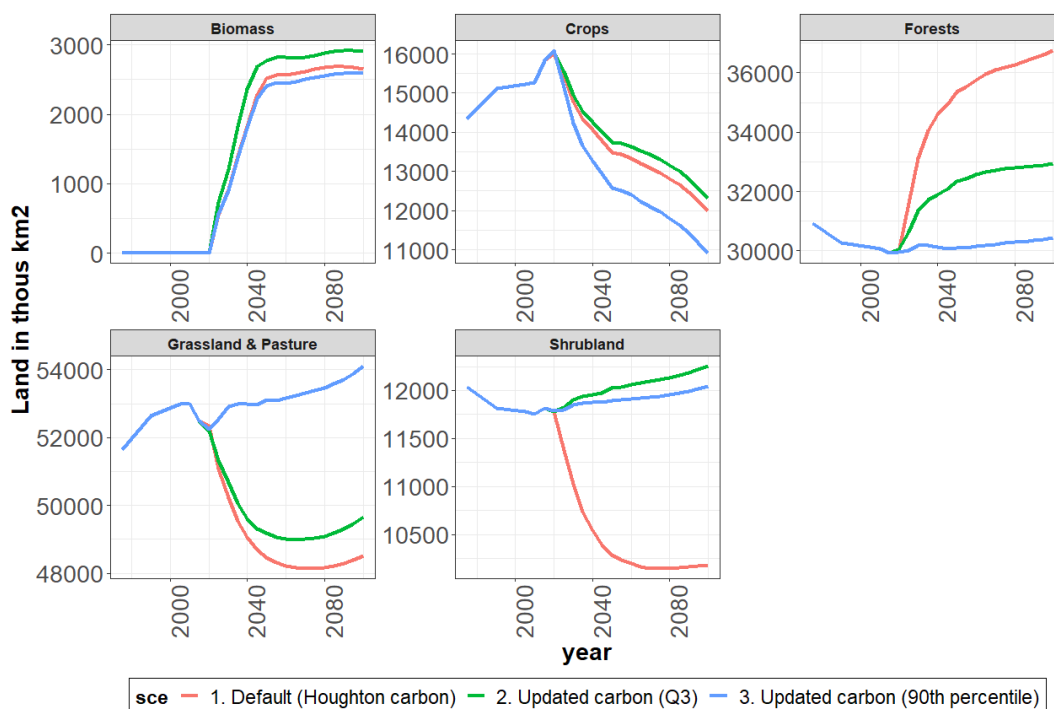


Figure 15: Global land allocation in GCAM under the SSP1 2p6 scenario by land type

The implementation of the spatially explicit carbon clearly improves land use responses and also suggests that high carbon sequestering shrubs can also be preserved as a part of nature based solutions to mitigate climate change. The robustness of these responses across other radiative forcing scenarios (implemented for more SSPs for example) and across other models need to be studied and is a subject worthy of exploration in a future paper.

5. Discussion and conclusion

In this paper we present a new dataset of grid cell level spatially explicit carbon harmonized with Moriai/GCAM land types. Our harmonized dataset presents carbon values for 3 pools (topsoil, above ground biomass and below ground biomass) for six statistical states for various land use types. Our dataset is available both at a 5 arcmin resolution and aggregated to 699 land regions. This dataset is specifically designed to enable initialization of spatially explicit carbon in IAMs and MSD models. This dataset was specifically designed to generate carbon values for GCAM, but can and should be extensible to other models. In the future, this dataset can be extended to include deeper soil (beyond 0-30 cms) so that land use responses in models can account for an additional deep soil carbon pool.



1
2 We noted that there are some limitations with respect to the carbon observations (both for soil
3 and vegetation) for the Tundra region. For example, we could find no data for 29% of the 5
4 arcmin gridcells for this land type. The biome mapping also needed to include several source
5 land types to enable an increase in data coverage for Tundra. This issue was likely caused by the
6 different definitions of Tundra land cover in different datasets. Recently, there have been efforts
7 dedicated to collecting carbon data specifically for this land type. These data should be
8 integrated in future releases of our data to address the current lack of data coverage.
9

10 As a part of our analysis, we observed that SoilGrids soil carbon values for cropland do not show
11 a depletion when compared to SoilGrids soil carbon values in unmanaged land. As discussed,
12 this is likely due the locations of sampling for cropland soil carbon. As a result, we reduced
13 Cropland soil carbon by 30% when we applied it to GCAM. If better/improved data on crop soil
14 carbon become available, our data should be updated with the same.
15

16 We have also noted that our current estimates of forest vegetation carbon are based on both
17 primary and secondary forests. This is due to the lack of availability of fine resolution (300 m)
18 land masks that distinguish between primary and secondary forests. As more data become
19 available related to forest cover types, a logical next step would be to break out different forest
20 types in our dataset.
21

22 Finally, our analysis showed that using the Q3 statistical state was most appropriate for GCAM
23 even though it resulted in an initialization of pre-industrial carbon value that was lower than
24 other estimates. Selection of the Q3 results in more accurate historical LUC emissions and the
25 model therefore spins up to a value that is close to other estimates in the literature in 2015. What
26 the correct initialization value is will differ from model to model and would require a similar
27 analysis.
28
29
30
31

32 **6. Data availability statement**

33 Final data are available for download here- <https://zenodo.org/record/7884615> (Narayan
34 et al., 2023) The data repository contains the following-

- 35 1. 72 rasters (4 land use types X 6 states X 3 carbon pools) at a 5 arcmin resolution
36 representative of carbon in 2010
- 37 2. 1 thematic raster which tracks 15 vegetation biomes for Unmanaged land use type (from 1.
38 above)
- 39 3. Tabular data file showing aggregated carbon densities for 6 states of carbon for 699 land
40 regions for soil (0-30cm), aboveground biomass and belowground biomass.
41
42

43 **7. Code availability statement**

44 As mentioned above, the data can be generated programmatically with scripts that are hosted on
45 GitHub (https://github.com/JGCRI/moirai/tree/master/ancillary/carbon_harmonization).



1 The process has been split into two steps where the computationally intensive stage 1
2 (approx.. 6 hours of processing) is optional with outputs made available in the repository.
3 The Stage 1 processing is performed using bash scripts which use the GDAL software
4 (Warmerdam, 2008). The second stage processing uses an R script and can be completed
5 for all carbon pools in approx. 15 minutes to generate the final 72 rasters and the final
6 tabular output file. We have also made available optional diagnostic functions in the R
7 script which can be used to validate results.

8
9 **Competing Interests declaration**- The authors have declared that none of the authors has
10 any competing interests.

11 **Acknowledgements**

12 This research was supported by the U.S. Department of Energy, Office of Science, as part
13 of research in Multi Sector Dynamics, Earth and Environmental System Modeling
14 Program. The Pacific Northwest National Laboratory is operated for DOE by Battelle
15 Memorial Institute under contract DE-AC05-76RL01830
16
17

18 **References**

- 19
20
21 Barnes, W. L., Xiong, X., & Salomonson, V. V. (2003). Status of terra MODIS and aqua MODIS. *Advances*
22 *in Space Research*, 32(11), 2099-2106.
23 Batjes, N. H., Ribeiro, E., Van Oostrum, A., Leenaars, J., Hengl, T., & Mendes de Jesus, J. (2017). WoSIS:
24 providing standardised soil profile data for the world. *Earth System Science Data*, 9(1), 1-14.
25 Calvin, K., Patel, P., Clarke, L., Asrar, G., Bond-Lamberty, B., Cui, R. Y., Di Vittorio, A., Dorheim, K.,
26 Edmonds, J., & Hartin, C. (2019). GCAM v5. 1: representing the linkages between energy, water,
27 land, climate, and economic systems. *Geoscientific Model Development*, 12(2), 677-698.
28 Cooper, H., Sjögersten, S., Lark, R., & Mooney, S. (2021). To till or not to till in a temperate ecosystem?
29 Implications for climate change mitigation. *Environmental Research Letters*, 16(5), 054022.
30 Di Vittorio, A. V., Vernon, C. R., & Shu, S. (2020). Moirai version 3: a data processing system to generate
31 recent historical land inputs for global modeling applications at various scales. *Journal of Open*
32 *Research Software*, 8(PNNL-SA-142149).
33 Erb, K.-H., Kastner, T., Plutzer, C., Bais, A. L. S., Carvalhais, N., Fetzel, T., Gingrich, S., Haberl, H., Lauk, C.,
34 & Niedertscheider, M. (2018). Unexpectedly large impact of forest management and grazing on
35 global vegetation biomass. *Nature*, 553(7686), 73-76.
36 Fang, Y., Liu, C., Huang, M., Li, H., & Leung, L. R. (2014). Steady state estimation of soil organic carbon
37 using satellite-derived canopy leaf area index. *Journal of Advances in Modeling Earth Systems*,
38 6(4), 1049-1064.
39 Friedlingstein, P., Jones, M. W., O'Sullivan, M., Andrew, R. M., Bakker, D. C., Hauck, J., Le Quéré, C.,
40 Peters, G. P., Peters, W., & Pongratz, J. (2022). Global carbon budget 2021. *Earth System Science*
41 *Data*, 14(4), 1917-2005.
42 Hengl, T., de Jesus, J. M., MacMillan, R. A., Batjes, N. H., Heuvelink, G. B., Ribeiro, E., Samuel-Rosa, A.,
43 Kempen, B., Leenaars, J. G., & Walsh, M. G. (2014). SoilGrids1km—global soil information based
44 on automated mapping. *PLoS one*, 9(8), e105992.
45 Houghton, R. (2005). Aboveground forest biomass and the global carbon balance. *Global Change*
46 *Biology*, 11(6), 945-958.



- 1 Houghton, R. A. (1999). The annual net flux of carbon to the atmosphere from changes in land use
2 1850–1990. *Tellus B*, 51(2), 298-313.
- 3 Hugelius, G., Strauss, J., Zubrzycki, S., Harden, J. W., Schuur, E., Ping, C.-L., Schirrmeyer, L., Grosse, G.,
4 Michaelson, G. J., & Koven, C. D. (2014). Estimated stocks of circumpolar permafrost carbon
5 with quantified uncertainty ranges and identified data gaps. *Biogeosciences*, 11(23), 6573-6593.
- 6 Jackson, R. B., Lajtha, K., Crow, S. E., Hugelius, G., Kramer, M. G., & Piñeiro, G. (2017). The ecology of soil
7 carbon: pools, vulnerabilities, and biotic and abiotic controls. *Annual Review of Ecology,
8 Evolution, and Systematics*, 48(1), 419-445.
- 9 Jungkunst, H. F., Göpel, J., Horvath, T., Ott, S., & Brunn, M. (2022). Global soil organic carbon–climate
10 interactions: Why scales matter. *Wiley Interdisciplinary Reviews: Climate Change*, e780.
- 11 Justice, C., Townshend, J., Vermote, E., Masuoka, E., Wolfe, R., Saleous, N., Roy, D., & Morisette, J.
12 (2002). An overview of MODIS Land data processing and product status. *Remote sensing of
13 Environment*, 83(1-2), 3-15.
- 14 Klein Goldewijk, K., Beusen, A., Doelman, J., & Stehfest, E. (2017). Anthropogenic land use estimates for
15 the Holocene–HYDE 3.2. *Earth System Science Data*, 9(2), 927-953.
- 16 Li, W., MacBean, N., Ciaia, P., Defourny, P., Lamarche, C., Bontemps, S., Houghton, R. A., & Peng, S.
17 (2018). Gross and net land cover changes in the main plant functional types derived from the
18 annual ESA CCI land cover maps (1992–2015). *Earth System Science Data*, 10(1), 219-234.
- 19 Liu, X., Yu, L., Si, Y., Zhang, C., Lu, H., Yu, C., & Gong, P. (2018). Identifying patterns and hotspots of
20 global land cover transitions using the ESA CCI Land Cover dataset. *Remote Sensing Letters*,
21 9(10), 972-981.
- 22 Meiyappan, P., & Jain, A. K. (2012). Three distinct global estimates of historical land-cover change and
23 land-use conversions for over 200 years. *Frontiers of earth science*, 6(2), 122-139.
- 24 Nachtergaele, F., van Velthuisen, H., Verelst, L., Batjes, N., Dijkshoorn, K., van Engelen, V., Fischer, G.,
25 Jones, A., & Montanarella, L. (2010). The harmonized world soil database. Proceedings of the
26 19th World Congress of Soil Science, Soil Solutions for a Changing World, Brisbane, Australia, 1-6
27 August 2010,
- 28 Kanishka B. Narayan, Alan Di Vittorio, Evan Margiotta, Seth A Spawn, & Holly Gibbs. (2023). Spatially
29 explicit re-harmonized terrestrial carbon densities for calibrating Integrated Multisectoral
30 Models (1.0.0) [Data set]. Zenodo. <https://doi.org/10.5281/zenodo.7884615>
- 31 Poggio, L., De Sousa, L. M., Batjes, N. H., Heuvelink, G., Kempen, B., Ribeiro, E., & Rossiter, D. (2021).
32 SoilGrids 2.0: producing soil information for the globe with quantified spatial uncertainty. *Soil*,
33 7(1), 217-240.
- 34 Pongratz, J., Schwingshackl, C., Bultan, S., Obermeier, W., Havermann, F., & Guo, S. (2021). Land use
35 effects on climate: current state, recent progress, and emerging topics. *Current Climate Change
36 Reports*, 1-22.
- 37 Popp, A., Calvin, K., Fujimori, S., Havlik, P., Humpenöder, F., Stehfest, E., Bodirsky, B. L., Dietrich, J. P.,
38 Doelmann, J. C., & Gusti, M. (2017). Land-use futures in the shared socio-economic pathways.
39 *Global Environmental Change*, 42, 331-345.
- 40 Ramankutty, N., & Foley, J. A. (1999). Estimating historical changes in land cover: North American
41 croplands from 1850 to 1992: GCTE/LUCC RESEARCH ARTICLE. *Global Ecology and
42 Biogeography*, 8(5), 381-396.
- 43 Sanderman, J., Hengl, T., & Fiske, G. J. (2017). Soil carbon debt of 12,000 years of human land use.
44 *Proceedings of the National Academy of Sciences*, 114(36), 9575-9580.
- 45 Scharlemann, J. P., Tanner, E. V., Hiederer, R., & Kapos, V. (2014). Global soil carbon: understanding and
46 managing the largest terrestrial carbon pool. *Carbon Management*, 5(1), 81-91.
- 47 Spawn, S. A., Sullivan, C. C., Lark, T. J., & Gibbs, H. K. (2020). Harmonized global maps of above and
48 belowground biomass carbon density in the year 2010. *Scientific Data*, 7(1), 1-22.



- 1 Thomson, A. M., Calvin, K. V., Chini, L. P., Hurtt, G., Edmonds, J. A., Bond-Lamberty, B., Frohling, S., Wise,
2 M. A., & Janetos, A. C. (2010). Climate mitigation and the future of tropical landscapes.
3 *Proceedings of the National Academy of Sciences*, *107*(46), 19633-19638.
- 4 Tifafi, M., Guenet, B., & Hatté, C. (2018). Large differences in global and regional total soil carbon stock
5 estimates based on SoilGrids, HWSD, and NCSCD: Intercomparison and evaluation based on field
6 data from USA, England, Wales, and France. *Global Biogeochemical Cycles*, *32*(1), 42-56.
- 7 van Asselen, S., & Verburg, P. H. (2012). AL and S system representation for global assessments and land-
8 use modeling. *Global Change Biology*, *18*(10), 3125-3148.
- 9 Walker, W. S., Gorelik, S. R., Cook-Patton, S. C., Baccini, A., Farina, M. K., Solvik, K. K., Ellis, P. W.,
10 Sanderman, J., Houghton, R. A., & Leavitt, S. M. (2022). The global potential for increased
11 storage of carbon on land. *Proceedings of the National Academy of Sciences*, *119*(23),
12 e2111312119.
- 13 Warmerdam, F. (2008). Open source approaches in spatial data handling. by Hall, GB & Leahy, MG
14 *Berlin, Heidelberg: Springer Berlin Heidelberg*, 87-104.
- 15 Wei, X., Shao, M., Gale, W., & Li, L. (2014). Global pattern of soil carbon losses due to the conversion of
16 forests to agricultural land. *Scientific reports*, *4*(1), 1-6.
- 17 Wieder, W. R., Boehnert, J., & Bonan, G. B. (2014). Evaluating soil biogeochemistry parameterizations in
18 Earth system models with observations. *Global Biogeochemical Cycles*, *28*(3), 211-222.
- 19 Wise, M., Calvin, K., Thomson, A., Clarke, L., Bond-Lamberty, B., Sands, R., Smith, S. J., Janetos, A., &
20 Edmonds, J. (2009). Implications of limiting CO₂ concentrations for land use and energy. *Science*,
21 *324*(5931), 1183-1186.
- 22