

# SELF-SUPERVISED MULTI-MODAL WORLD MODEL WITH 4D SPACE-TIME EMBEDDING

**Lance Legel\***

Ecological Intelligence Lab  
Ecodash.ai

**Qin Huang**

School of Complex Adaptive Systems  
Arizona State University

**Brandon Voelker**

Geosensing Systems Engineering & Sciences Lab  
University of Houston

**Daniel Neamati**

Navigation & Autonomous Vehicles Lab  
Stanford University

**Patrick Alan Johnson**

Earth System Lab  
Allen Institute for Artificial Intelligence

**Favyen Bastani**

Earth System Lab  
Allen Institute for Artificial Intelligence

**Jeff Rose**

Spatial Intelligence Lab  
SpatialLogic.com

**James Ryan Hennessy**

Department of Computer Science  
Georgia Tech Institute of Technology

**Robert Guralnick**

Florida Museum of Natural History  
University of Florida

**Douglas Soltis**

Florida Museum of Natural History  
University of Florida

**Pamela Soltis**

Florida Museum of Natural History  
University of Florida

**Shaowen Wang**

NSF Institute for Geospatial Understanding  
University of Illinois Urbana-Champaign

## ABSTRACT

We present *DeepEarth*, a self-supervised multi-modal world model with *Earth4D*, a novel planetary-scale 4D space-time positional encoder. *Earth4D* extends 3D multi-resolution hash encoding to include time, efficiently scaling across the planet over centuries with sub-meter, sub-second precision. Multi-modal encoders (*e.g.* vision-language models) are fused with *Earth4D* embeddings and trained via masked reconstruction. We demonstrate *Earth4D*'s expressive power by achieving state-of-the-art performance on an ecological forecasting benchmark. *Earth4D* with learnable hash probing surpasses a multi-modal foundation model pre-trained on substantially more data. Access open source code and download models at: <https://github.com/legel/deepearth>.

## 1 DEEPEARTH ARCHITECTURE

*DeepEarth* is a self-supervised multi-modal world model that learns unified representations of Earth observation data across space and time. As seen in Figure 1, the architecture processes multi-modal inputs (*e.g.* vision, language, sensor data) sampled around spatio-temporal events. The *Earth4D* encoder maps continuous space-time coordinates (*latitude*, *longitude*, *elevation*, *time*) to learnable positional embeddings, which are fused with embeddings from modality-specific encoders and processed as tokens in an autoencoder context window. Inspired by PerceiverIO (Jaegle et al., 2021), V-JEPA 2 (Assran et al., 2025), Galileo (Tseng et al., 2025), and AlphaEarth (Brown et al., 2025), *DeepEarth* learns to generatively reconstruct and simulate joint distributions of multi-modal data.

---

\*Correspondence: lance@ecodash.ai

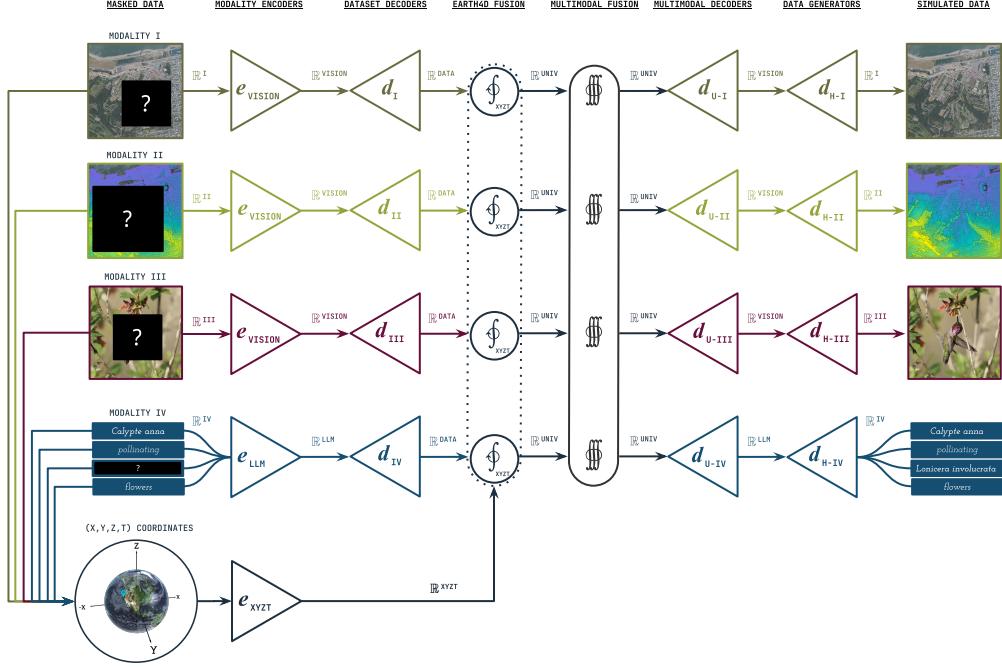


Figure 1: **DeepEarth Overview.** Masked multi-modal data (e.g. images, text) sampled around an event (e.g. pollination) are encoded and fused with Earth4D space-time embeddings. These universal tokens are jointly encoded, and then masked data is inductively decoded and simulated.

## 2 EARTH4D ARCHITECTURE

Following Grid4D (Jiawei et al., 2024), Earth4D extends NVIDIA’s multi-resolution hash encoding (Müller et al., 2022) to four dimensions (Figure 2) by concatenating features from one spatial ( $xyz$ ) and three spatio-temporal ( $xyt$ ,  $yzt$ ,  $xzt$ ) grids. Implemented as a standalone PyTorch module with massively parallelizable CUDA kernels, Earth4D is suitable for integration into larger models.

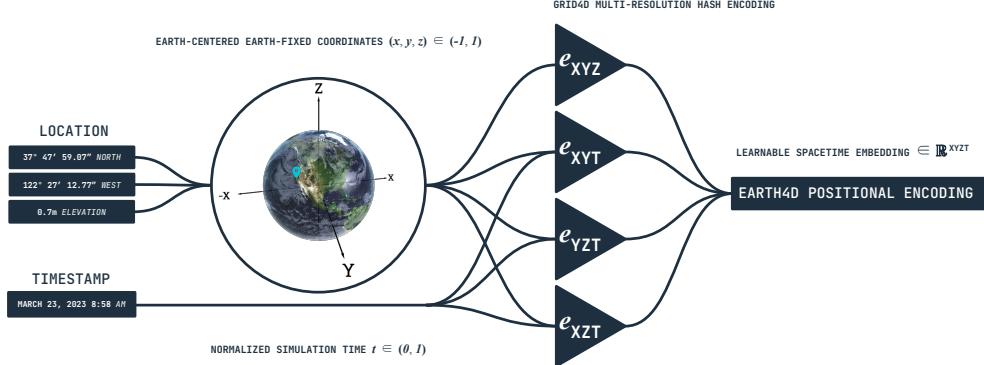
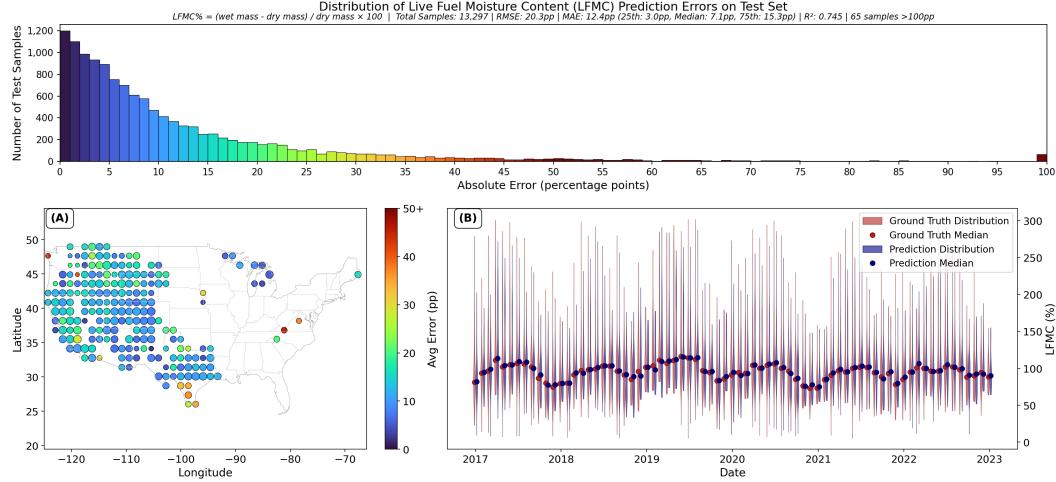


Figure 2: **Earth4D Space-Time Positional Encoding.** A planetary-scale 4D encoder with fully decomposable spatio-temporal representation. Four grids ( $xyz$ ,  $xyt$ ,  $yzt$ ,  $xzt$ ) are each learned in 3D space and computed in parallel. Each grid has multiple resolution levels (Appendix A), enabling deep learning of complex joint distributions in multi-modal data across space-time scales.



**Figure 3: Earth4D LFMC Prediction Performance.** (Top) Distribution of absolute errors in percentage point predictions across 13,297 test samples, showing median error of 7.1pp. (A) Geographic error distribution across CONUS shows low error in well-sampled regions. (B) Temporal predictions closely track ground truth LFMC measurements across seasons (2017–2023).

Earth4D’s hash encoding compresses spatial features into a fixed memory budget, but different coordinates can map to the same memory location (collisions). We integrate learned hash probing (Takikawa et al., 2023), an end-to-end differentiable system that learns optimal memory allocation patterns for the data. This yields substantial performance improvements across tasks (Appendix B).

### 3 EARTH4D EXPERIMENTAL VALIDATION

#### 3.1 LIVE FUEL MOISTURE CONTENT PREDICTION

**Dataset.** Live Fuel Moisture Content (LFMC) measures the percentage of water in vegetation relative to its dry weight, a critical indicator for wildfire risk assessment. We evaluate Earth4D on Globe-LFMC 2.0 (Yebra et al., 2024), a global ecological forecasting benchmark containing field measurements across diverse plant species, geographic regions, and temporal periods.

**Baseline Model.** We compare against Galileo (Johnson et al., 2025; Tseng et al., 2025), a pre-trained Vision Transformer processing multi-modal remote sensing data (Appendix C).

**Architecture.** Earth4D encodes  $(x,y,z,t)$  into a 192D vector, concatenated with a learnable species embedding initialized randomly (no prior knowledge). An MLP then predicts LFMC %.

**Results.** Earth4D achieves MAE 12.1pp and  $R^2$  0.755, surpassing Galileo (MAE 12.6pp,  $R^2$  0.72) using only  $(x,y,z,t)$  coordinates and species embeddings (Table 1).

Model	Data Inputs	MAE (pp)	RMSE (pp)	$R^2$
Galileo (Pre-Trained)	$(x,y,z,t)$ + Species Type + Remote Sensing	12.6	<b>18.9</b>	0.72
Earth4D (Learned Hashing)	$(x,y,z,t)$ + Species Name	<b>12.1</b>	19.9	<b>0.755</b>

**Table 1: State-of-the-Art Ecological Forecasting Benchmark.** Earth4D surpasses the pre-trained Galileo foundation model without satellite imagery, weather data, or topography.

#### 3.2 RGB AERIAL IMAGERY RECONSTRUCTION

We evaluate Earth4D’s ability to infer RGB pixels from  $(x,y,z,t)$  inputs with objective  $(x,y,z,t) \rightarrow (r,g,b)$ . Using USGS 3DEP LiDAR (Stoker & Miller, 2022; Sugabaker et al., 2014) and USDA NAIP imagery (USDA, 2003–present) paired by Allred et al. (2025), we train on 5.8M coordinate-color pairs from Houston coastal wetlands (Figure 4).



Figure 4: **RGB Reconstruction from LiDAR Elevation.** Houston coastal wetlands, 2018. *Left to right:* LiDAR height, ground truth, baseline, learned probing (18% lower loss).

## REFERENCES

- John T Abatzoglou, Solomon Z Dobrowski, Sean A Parks, and Katherine C Hegewisch. TerraClimate, a high-resolution global dataset of monthly climate and climatic water balance from 1958–2015. *Scientific Data*, 5:1–12, 2018.
- Brady W Allred, Sarah E McCord, and Scott L Morford. Canopy height model and NAIP imagery pairs across CONUS. *Scientific Data*, 12(1):322, 2025.
- Mido Assran, Adrien Bardes, David Fan, Quentin Garrido, Russell Howes, Mojtaba, Komeili, Matthew Muckley, Ammar Rizvi, Claire Roberts, Koustuv Sinha, Artem Zhulus, Sergio Arnaud, Abha Gejji, Ada Martin, Francois Robert Hogan, Daniel Dugas, Piotr Bojanowski, Vasil Khalidov, Patrick Labatut, Francisco Massa, Marc Szafraniec, Kapil Krishnakumar, Yong Li, Xiaodong Ma, Sarath Chandar, Franziska Meier, Yann LeCun, Michael Rabbat, and Nicolas Ballas. V-JEPA 2: Self-supervised video models enable understanding, prediction and planning, 2025. URL <https://arxiv.org/abs/2506.09985>.
- Christopher F. Brown, Michal R. Kazmierski, Valerie J. Pasquarella, William J. Rucklidge, Masha Samsikova, Chenhui Zhang, Evan Shelhamer, Estefania Lahera, Olivia Wiles, Simon Ilyushchenko, Noel Gorelick, Lihui Lydia Zhang, Sophia Alj, Emily Schechter, Sean Askay, Oliver Guinan, Rebecca Moore, Alexis Boukouvalas, and Pushmeet Kohli. AlphaEarth Foundations: An embedding field model for accurate and efficient global mapping from sparse label data, 2025. URL <https://arxiv.org/abs/2507.22291>.
- Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, Jakob Uszkoreit, and Neil Houlsby. An image is worth 16x16 words: Transformers for image recognition at scale. In *International Conference on Learning Representations*, 2021. URL <https://openreview.net/forum?id=YicbFdNTTy>.
- Matthias Drusch, Umberto Del Bello, Sébastien Carlier, Olivier Colin, Veronica Fernandez, Ferran Gascon, Bianca Hoersch, Claudia Isola, Paolo Laberinti, Philippe Martimort, et al. Sentinel-2: ESA’s optical high-resolution mission for GMES operational services. *Remote Sensing of Environment*, 120:25–36, 2012.
- Tom G Farr and Mike Kobrick. Shuttle radar topography mission produces a wealth of data. *Eos, Transactions American Geophysical Union*, 81:583–585, 2000.
- Andrew Jaegle, Sebastian Borgeaud, Jean-Baptiste Alayrac, Carl Doersch, Catalin Ionescu, David Ding, Skanda Koppula, Daniel Zoran, Andrew Brock, Evan Shelhamer, et al. Perceiver IO: A general architecture for structured inputs & outputs. In *International Conference on Learning Representations*, 2021.
- Xu Jiawei, Fan Zexin, Yang Jian, and Xie Jin. Grid4D: 4D decomposed hash encoding for high-fidelity dynamic gaussian splatting. *The Thirty-eighth Annual Conference on Neural Information Processing Systems*, 2024.
- Patrick Alan Johnson, Gabriel Tseng, Yawen Zhang, Heather Heward, Virginia Sjahli, Favyen Bastani, Joseph Redmon, and Patrick Beukema. High-resolution live fuel moisture content (LFMC) maps for wildfire risk from multimodal earth observation data, 2025. URL <https://arxiv.org/abs/2506.20132>.
- Joaquín Muñoz Sabater. ERA5-Land monthly averaged data from 1950 to present. Copernicus Climate Change Service (C3S) Climate Data Store (CDS), 2019.
- Thomas Müller, Alex Evans, Christoph Schied, and Alexander Keller. Instant neural graphics primitives with a multiresolution hash encoding. *ACM Transactions on Graphics*, 41(4):1–15, July 2022. ISSN 1557-7368. doi: 10.1145/3528223.3530127. URL <http://dx.doi.org/10.1145/3528223.3530127>.
- Jason Stoker and Barry Miller. The accuracy and consistency of 3D elevation program data: A systematic analysis. *Remote Sensing*, 14(4):940, 2022.

Larry J Sugarbaker, Eric W Constance, Hans Karl Heidemann, Allyson L Jason, Vicki Lukas, David L Saghy, and Jason M Stoker. The 3D elevation program initiative: a call for action. Technical report, US Geological Survey, 2014.

Towaki Takikawa, Thomas Müller, Merlin Nimier-David, Alex Evans, Sanja Fidler, Alec Jacobson, and Alexander Keller. Compact neural graphics primitives with learned hash probing, 2023. URL <https://arxiv.org/abs/2312.17241>.

Ramon Torres, Paul Snoeij, Dirk Geudtner, David Bibby, Malcolm Davidson, Evert Attema, Pierre Potin, Börn Rommen, Nicolas Floury, Mike Brown, et al. GMES Sentinel-1 mission. *Remote Sensing of Environment*, 120:9–24, 2012.

Gabriel Tseng, Anthony Fuller, Marlena Reil, Henry Herzog, Patrick Beukema, Favyen Bastani, James R Green, Evan Shelhamer, Hannah Kerner, and David Rolnick. Galileo: Learning global & local features of many remote sensing modalities. In *Forty-second International Conference on Machine Learning*, 2025.

USDA. National Agriculture Imagery Program (NAIP). <https://naip-usdaonline.hub.arcgis.com/>, 2003–present. Accessed: November 23, 2025.

Marta Yebra, Gianluca Scortechini, Karine Adeline, Nursema Aktepe, Turkia Almoustafa, Avi Bar-Massada, María Eugenia Beget, Matthias Boer, Ross Bradstock, Tegan Brown, et al. Globe-LFMC 2.0, an enhanced and updated dataset for live fuel moisture content research. *Scientific Data*, 11(1):332, 2024.

## APPENDICES

## A EARTH4D RESOLUTION SPECIFICATIONS

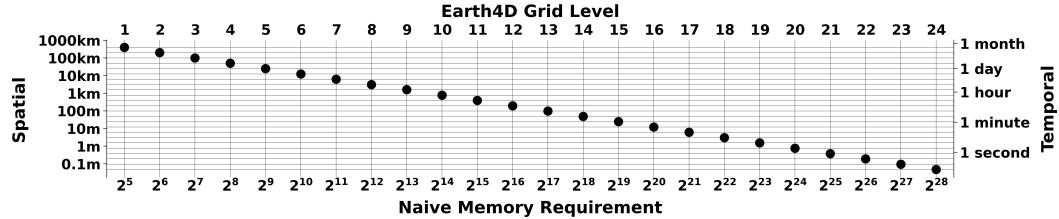


Figure 5: **Earth4D Space-Time Scales.** Default  $24 \times 24 \times 24$  levels for each  $xyz$ ,  $xyt$ ,  $yzt$ ,  $xzt$  grid. Each level stores up to  $2^{22}$  entries, with each entry storing a 2D feature. Requires 724M trainable parameters ( $\sim 11$  GB GPU memory during training). Parallelizable across levels and spatio-temporal boundaries. Outputs 192D per  $(x, y, z, t)$  coordinate from 4 grids  $\times$  24 levels  $\times$  2D feature per level. Hashing saves memory vs. naive requirement, e.g.,  $(2^{28})^3 = 10^{25}$  at level 24.

## B LEARNED HASH PROBING AND ABLATION STUDIES

### B.1 HASH COLLISION SIMULATIONS

Scenario	Spatial	Temporal	Description
<i>Uniform Random</i>	Global	Full	Uniform Earth surface sampling
<i>Continental Sparse</i>	North America	Full	Sparse continental coverage
<i>Moderate Spatial Cluster</i>	10km × 10km	Full	City-scale clustering
<i>Moderate Temporal Cluster</i>	1k locations	Distributed	Temporal sampling at fixed locations
<i>Moderate Spatiotemporal</i>	1km × 1km	1 hour	Neighborhood-scale event
<i>Extreme Spatial Single</i>	10m × 10m	Full	Building-scale dense clustering
<i>Extreme Spatial Multi</i>	10 × (10m × 10m)	Full	10 dense clusters worldwide
<i>Extreme Temporal Single</i>	Global	1 hour	Global snapshot
<i>Extreme Temporal Multi</i>	Global	10 × (1 hour)	10 temporal snapshots
<i>Time Series</i>	10k locations	100 steps	Regular temporal sampling

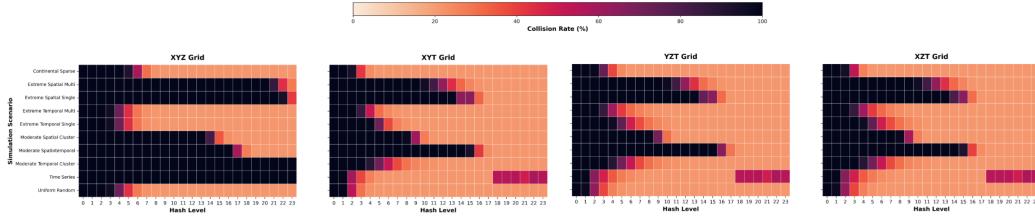


Figure 6: **Earth4D Hash Collision Analysis.** (Table) 10 ( $x, y, z, t$ ) point distribution scenarios that were simulated to analyze hash collisions in Earth4D memory. (Graph) Shows results for 1M point simulations across all 24 levels.

### B.2 PERFORMANCE IMPROVEMENTS

Standard multi-resolution hash encoding without learned probing obtains RMSE 26.0pp, MAE 16.6pp, and  $R^2$  0.58 (800M parameters,  $2^{22}$  hash capacity). Integrating learned hash probing (Takikawa et al., 2023), which learns to select optimal hash table indices from a candidate set, yields RMSE 19.9pp, MAE 12.1pp, and  $R^2$  0.755—a 27.1% MAE reduction and 30.2%  $R^2$  improvement. Extreme compression to 5M parameters (99.3% reduction,  $2^{14}$  hash capacity) achieves MAE 15.0pp/ $R^2$  0.668, outperforming the 800M baseline by 14.7% in  $R^2$  with 4× training speedup and 93% memory reduction. On RGB reconstruction, learned probing reduces validation loss by 18%. These gains result from collision reduction (33% at 1M points) and learned shared features across memory indices, allowing the model to discover meaningful spatio-temporal patterns.

## C BENCHMARK SPECIFICATIONS

### C.1 GALILEO BASELINE MODEL

Galileo (Johnson et al., 2025; Tseng et al., 2025) is a Vision Transformer (Dosovitskiy et al., 2021) (5.3M parameters) pre-trained by the Allen Institute for AI. It processes Sentinel-2 optical imagery (Drusch et al., 2012), Sentinel-1 SAR (Torres et al., 2012), VIIRS night lights, ERA-5 weather (Muñoz Sabater, 2019), TerraClimate soil/water data (Abatzoglou et al., 2018), SRTM topography (Farr & Kobrick, 2000),  $(x,y,z,t)$  coordinates, and species type. We use the Allen Institute for AI's exact Globe-LFMC 2.0 (Yebra et al., 2024) train/test split (76,467/13,297) to directly compare against this benchmark.