# DynamicEarthNet: towards FedSPD

April 1, 2025

# Overview

# DynamicEarthNet dataset

See [TKW+22].

Multi-spectral images

- ▶ 75 Areas of Interest (AOI)
- ▶ Satellite images collection: daily $2018/01/01 \rightarrow 2019/12/31$
- ▶ 730 images/AOI $\rightarrow 54,750$ satellite images in total
- ▶ 4 channels: RGB+NIR
- ▶ Spatial resolution: 3 meters
- ▶ Image size: $1024 \times 1024$

One time-series: $x \in \mathbb{R}^{T \times H \times W \times 4}$ with $T$:number of days in the month, $H = W = 1024$.

- Labels: $24 \times 75 = 1,800$
- 7 different (imbalanced) classes:

| class name | % | #AOIs | color |
|---|---|---|---|
| impervious surface | 7.1 | 70 | |
| agriculture | 10.3 | 37 | |
| forest & other vegetation | 44.9 | 71 | |
| wetlands | 0.7 | 24 | |
| soil | 28.0 | 75 | |
| water | 8.0 | 58 | |
| snow & ice | 1.0 | 2 | |

No single AOI contains all 7 classes.
For a given time series, label:

$$y \in \{0, \ldots, 6\}^{T \times H \times W}$$

(a)                                    (b)

Figure: Two random images from different time series id (and different directories)

**Challenges:**

Issue 1 Voluminous data (hundreds of GB)

Issue 2 Slow to move all images to a single processing server

Solution Issue 1 client $\rightarrow$ smaller dataset & own stat. heterogeneity with clear interpretation

Solution Issue 2 No exchange of raw datasets, just each client's updates

▶ Each AOI $\rightarrow$ local phenomena

▶ FL: adapt to local distrib and aggregate knowledge for better generalization

  ⋆ Idea 1 one AOI $=$ one client $\rightarrow$ too many clients (75);

  ⋆ Idea 2 create clients based on dominant LULC categories across AOI

    1. Urban areas
    2. Agricultural areas
    3. Forest-dominated areas
    4. Water-dominated areas
    5. Mixed land-use areas

## Framework

$K$ clients. Each client $k$ has dataset $D^{(k)}$ of $n_k$ pixel annotation pairs $\{\mathbf{X}_i^{(k)}, y_i^{(k)}\}_{i=1,\dots,n_k}$, where

- $\mathbf{X}_i^{(k)}$ is a covariance matrix in $\mathbb{R}^{4T \times 4T}$ (with $T = 31$),
- $y_i^{(k)} \in \{0, \dots, 6\}$ is a semantic mask

Total number of samples $n := \sum_{k=1}^{K} n_k$

Loss function $\ell$ (CE)

Forward pass: SPD-Net $f_\Theta(\cdot)$ (see [HVG17, BSB$^+$19])

**SPD Feature Extraction**



Each client's SPD-net with $L$ layers has learnable parameters

$$\Theta = (\underbrace{W_1, \ldots, W_L}_{BiMap}, \underbrace{G_1, \ldots, G_L}_{BN}, \underbrace{\xi_1, \ldots, \xi_d}_{FC}) \in \mathcal{M}$$

with

$$\mathcal{M} := \left( \prod_{l=1}^{L} \mathrm{St}(d_l, d_{l+1}) \right) \times \left( \prod_{l=1}^{L} \mathrm{SPD}_{d_l} \right) \times \mathbb{R}^d.$$

Local objective functions For each client $k = 1, \ldots, K$,

$$F_k(\Theta) := \frac{1}{n_k} \sum_{i=1}^{n_k} \ell\big(f_\Theta(x_i^{(k)}), y_i^{(k)}\big).$$

Global objective function

$$\arg\min_{\Theta \in \mathcal{M}} F(\Theta) := \sum_{k=1}^{K} p_k F_k(\Theta), \quad p_k = \frac{n_k}{n}.$$

# Federated procedure

See [LM22].

## Algorithm

$\star$ **Initialization** $\Theta_0$ to all clients (random)

$\star$ **Outer loop** For $t = 0, \ldots, T - 1$ (rounds):

1. A subset $\mathcal{S}_t$ of $r \leq K$ clients is sampled (random)
2. Server broadcasts $\Theta_t$ to clients.
3. Each client $k \in \mathcal{S}_t$ runs $E_k$ local updates on $\Theta$ with initialization $\Theta_t$.
4. Aggregation of each client's local updated parameters via barycenters $\rightarrow \Theta_{t+1}$

$\star$ **Global output** Final aggregated param. $\Theta_T$ or $\mathrm{Unif}(\{\Theta_0, \ldots, \Theta_T\})$

If each local function $F_k$ is $L_k$-smooth and geodesically convex, then the algorithm converges in $O(1/T)$ (provided some curvature/manifold additional assumptions).

Let $\theta$ denote a Riemannian parameter (either one of the $W_l$ or $G_l$).

- Initialize $\theta_0^{(k)} = \theta_t$
- For $s = 1, \ldots, E_k$:

$$\theta_{s+1}^{(k)} \;=\; R_{\theta_s^{(k)}}\Big(-\eta^{(k)}\, g^{(k)}\Big),$$

  where the local Riemannian gradient is

$$g^{(k)} \;=\; \operatorname{grad} F_k\big(\theta_s^{(k)}\big) \;-\; \mathcal{T}_{\theta_t \to \theta_s^{(k)}}\Big(\operatorname{grad} F_k(\theta_t) \;-\; \operatorname{grad} F(\theta_t)\Big),$$

  - The blue term is the optional **SVR correction** (variance reduction).

# Differential Privacy as a model variant

Besides the variance-reduction approach, we can also inject **DP noise** (see [HHJM24a, HMJG24]) on each local gradient to protect sensitive data.

**Key idea:** Add tangent-space Gaussian noise before retraction/exponential update.

- ▶ **DP** can be combined with or without the SVR term.
- ▶ Each client $k$ ensures local DP by perturbing its gradient:

$$g_{\text{DP}}^{(k)} = \text{clip}_\tau(g^{(k)}) + \varepsilon_t, \quad \text{where} \quad \varepsilon_t \sim \mathcal{N}_{\theta_s^{(k)}}(0, \sigma^2).$$

$\mathcal{N}_{\theta_s^{(k)}}(0, \sigma^2)$ is *tangent-space Gaussian* distribution at $\theta_s^{(k)}$ and $\text{clip}_\tau : T_\theta \mathcal{M} \to T_\theta \mathcal{M}$ is such that

$$\text{clip}_\tau(\xi) = \min\left\{\frac{\tau}{\|\xi\|_\theta}, 1\right\}.$$

**Privacy Guarantee:** With an appropriate choice of $\sigma$ and composition rules, the entire FL pipeline can provide an $(\varepsilon, \delta)$-DP guarantee.

# Riemannian local updates with optional DP

- Initialize $\theta_0^{(k)} = \theta_t$
- For $s = 1, \ldots, E_k$: over a selected minibatch,

$$g^{(k)} = \operatorname{grad} F_k(\theta_s^{(k)}) - \mathcal{T}_{\theta_t \to \theta_s^{(k)}}\Big(\operatorname{grad} F_k(\theta_t) - \operatorname{grad} F(\theta_t)\Big),$$

$$g_{\mathrm{DP}}^{(k)} = \operatorname{clip}_\tau(g^{(k)}) + \varepsilon_s, \qquad \varepsilon_s \sim \mathcal{N}_{\theta_s^{(k)}}(0, \sigma^2).$$

Then update via retraction (or exponential map):

$$\theta_{s+1}^{(k)} = R_{\theta_s^{(k)}}\Big(-\eta^{(k)} g_{\mathrm{DP}}^{(k)}\Big).$$

- **Red** term is the **DP noise** addition (optional)
- **Blue** term is still the **SVR correction** (optional).

Let $\xi \in \mathbb{R}^{d_{\mathrm{FC}}}$ be the parameter vector for the final fully-connected layer. During each local epoch on client $k$:

$$\xi_0^{(k)} = \xi_t, \quad \text{and for } s = 0, \ldots, E_k - 1:$$

$$\xi_{s+1}^{(k)} = \xi_s^{(k)} - \eta^{(k)} \left( \nabla F_k\big(\xi_s^{(k)}\big) - \big[\nabla F_k(\xi_t) - \nabla F(\xi_t)\big] + \varepsilon_s \right),$$

where:

- $\nabla F_k(\xi_s^{(k)})$: local Euclidean gradient on client $k$.
- $\big[\nabla F_k(\xi_t) - \nabla F(\xi_t)\big]$ is the **SVR correction** term.
- $\varepsilon_s \sim \mathcal{N}(0, \sigma^2 I)$ is the **DP noise** term.
- After $E_k$ steps, client $k$ sends $\xi_{E_k}^{(k)}$ to the server.

After each client $k \in \mathcal{S}_r$ sends back $\theta_{E_k}^{(k)}$, we compute the Riemannian (Fréchet) mean on the manifold:

$$\theta_{t+1} = \arg \min_{\theta \in \mathcal{M}} \sum_{k \in \mathcal{S}_r} p_k \, d\big(\theta, \, \theta_{E_k}^{(k)}\big)^2.$$

▶ In practice, we use a gradient-descent-based procedure (Karcher flow):

$$\theta_{t+1} \leftarrow \operatorname{Exp}_{\theta_t}\!\Big(-\gamma \sum_{k \in \mathcal{S}_r} p_k \operatorname{Log}_{\theta_t}(\theta_{E_k}^{(k)})\Big),$$

where $\gamma$ is a step size for the Karcher flow.

See [BLSB25].

For the BiMap parameters $W_l$ living on the Stiefel manifold, we can use:

▶ **Karcher flow** or Riemannian gradient-based aggregator:

$$W_{t+1} = \arg\min_{W \in \mathrm{St}(d_l, d_{l+1})} \sum_{k \in S_t} p_k \, d^2\big(W, \, W_\tau^{(k)}\big),$$

which can be iteratively approximated (exp/log maps).

▶ **Projection-based aggregator**:

$$W_{t+1} = \mathrm{uf}\Big(\underbrace{\sum_{k \in S_t} p_k \, W_\tau^{(k)}}_{\text{arithmetic mean}}\Big) \quad \text{or} \quad W_{t+1} = \mathrm{qf}\Big(\sum_{k \in S_t} p_k \, W_\tau^{(k)}\Big).$$

No iterative scheme needed, computationally cheaper (simple retractions/liftings).

At round $t$, we receive each client's updated $\xi_{E_k}^{(k)}$ for $k \in \mathcal{S}_t$. To get $\xi_{t+1}$, standard weighted average (FedAvg):

$$\xi_{t+1} \ = \ \sum_{k \in \mathcal{S}_t} \frac{n_k}{\sum_{j \in \mathcal{S}_t} n_j} \, \xi_{E_k}^{(k)}$$

No need for any Riemannian retraction or manifold-based aggregator.

## Average gradient stream approach

See [HHJM24b].

Key ideas :

1. Aggregate each client's mini-batch gradient with vector transport to server.
2. Update global model.

Pros :

- ▶ No communication of parameters $\rightarrow$ no expensive exp/log maps.
- ▶ Works for general manifolds (not just embedded or compact manifolds).
- ▶ Sublinear convergence for non-convex objectives (fixed step-size).

Cons :

- ▶ High variance of aggregating gradients, cost of vector transport
- ▶ Memory overhead on clients (store and transmit cumulative transported gradients)

# RFedAGS algorithm

⋆ **Outer loop** For $t = 0, \ldots, T - 1$, server broadcasts $x_t$.

Each client $k = 1, \ldots, K$ does $\theta_0^{(k)} \leftarrow x_t$ ; $g_0^{(k)} \leftarrow 0$.

1. For $s = 0, \ldots, E_k$:
   - Sample mini-batch $B_s^{(k)}$, compute local gradient

$$\eta_s^{(j)} \leftarrow -\alpha_t^{(k)} \frac{1}{\# B_s^{(k)}} \sum_{s \in B_s^{(k)}} \mathsf{clip}_\tau \mathsf{grad} f(\theta_s^{(k)}) + \varepsilon_s^{(k)}$$

2. Retraction-based update: $\theta_{s+1}^{(k)} = R_{\theta_s^{(k)}}(\eta_{t,s}^{(k)})$
3. Transport gradient to server's tangent space and accumulate:

$$g_{s+1}^{(k)} = g_s^{(k)} + \mathcal{T}_{x_{t,s}^{(k)} \to x_t}(\eta_s^{(k)}).$$

⋆ **Aggregation** : Each client uploads $g_{E_k}^{(k)}$ and server update:

$$x_{t+1} \leftarrow R_{x_t}\left( \frac{1}{K} \sum_{k=1}^{K} g_{E_k}^{(k)} \right).$$

⋆ **Output** $x_T$. *In red, optional DP-RSGD style.*

# References

Florent Bouchard, Nils Laurent, Salem Said, and Nicolas Le Bihan.
Beyond r-barycenters: an effective averaging method on stiefel and grassmann manifolds.
*arXiv preprint arXiv:2501.11555*, 2025.

Daniel Brooks, Olivier Schwander, Frédéric Barbaresco, Jean-Yves Schneider, and Matthieu Cord.
Riemannian batch normalization for spd neural networks.
*Advances in Neural Information Processing Systems*, 32, 2019.

Zhenwei Huang, Wen Huang, Pratik Jawanpuria, and Bamdev Mishra.
Federated learning on riemannian manifolds with differential privacy.
*arXiv preprint arXiv:2404.10029*, 2024.

Zhenwei Huang, Wen Huang, Pratik Jawanpuria, and Bamdev Mishra.
Riemannian federated learning via averaging gradient stream.
*arXiv preprint arXiv:2409.07223*, 2024.

Andi Han, Bamdev Mishra, Pratik Jawanpuria, and Junbin Gao.
Differentially private riemannian optimization.
*Machine Learning*, 113(3):1133–1161, 2024.

Zhiwu Huang and Luc Van Gool.
A riemannian network for spd matrix learning.
In *Proceedings of the AAAI conference on artificial intelligence*, volume 31, 2017.

Jiaxiang Li and Shiqian Ma.
Federated learning on riemannian manifolds.
*arXiv preprint arXiv:2206.05668*, 2022.

Aysim Toker, Lukas Kondmann, Mark Weber, Marvin Eisenberger, Andrés Camero, Jingliang Hu, Ariadna Pregel Hoderlein, Çağlar Şenaras,
Timothy Davis, Daniel Cremers, et al.
Dynamicearthnet: Daily multi-spectral satellite dataset for semantic change segmentation.