

Multi-tasking, Meta-learning and Federated Learning

From Improved Performance to Privacy Preservation

Anshul Thakur

Department of Engineering Science
University of Oxford

7 Nov, 2025

Table of Contents

Preliminary Concepts

Transforms

Loss Landscape & Mode Connections

Multi-tasking

Adaptive Parameter Optimisation

Task Groupings

Gradient-based Meta-Learning

Federated Learning

Data Democratisation & Healthcare

Data Encoding and Condensation

Uncertainty in LLMs



Table of Contents

Preliminary Concepts

Transforms

Loss Landscape & Mode Connections

Multi-tasking

Adaptive Parameter Optimisation

Task Groupings

Gradient-based Meta-Learning

Federated Learning

Data Democratisation & Healthcare

Data Encoding and Condensation

Uncertainty in LLMs



What are Transforms?



- ▶ **Transform** data from **one** domain to **other** domain
 - ▶ E.g. Time domain to **frequency** domain
- ▶ Why transform data?
 - ▶ **Analysis** of data may be **easier** in transformed domain
 - ▶ Certain operations may only be **feasible** in transformed domain
 - ▶ **Privacy Preservation!** (potentially)
- ▶ **Linear Transformation** functions are **invertible**
- ▶ **Deep neural networks** are **non-linear non-invertible** transforms!!!

Transforms around us!



- ▶ **Discrete Cosine Transform:** Represents a finite signal as a **combination of cosine functions** of different frequencies
- ▶ Allows **compression** of signals. Used in:
 - ▶ digital images (**JPEG**)
 - ▶ digital video (**MPEG**)
 - ▶ digital audio (**MP3, Dolby Digital**)
 - ▶ speech coding
- ▶ Made **streaming** music, games and Netflix possible

Combination of Vectors



- ▶ Let $\{\mathbf{v}_1, \mathbf{v}_2, \mathbf{v}_3 \dots \mathbf{v}_n\}$ be a set of N vectors and $\{\alpha_1, \alpha_2, \alpha_3 \dots \alpha_n\}$ be a set of N scalars

- ▶ **Linear combination:** $\mathbf{y} = \sum_{i=1}^N \alpha_i \mathbf{v}_i$

- ▶ **Affine combination:** $\mathbf{y} = \sum_{i=1}^N \alpha_i \mathbf{v}_i$ where $\sum_{i=1}^N \alpha_i = 1$

- ▶ **Conic combination:** $\mathbf{y} = \sum_{i=1}^N \alpha_i \mathbf{v}_i$ where $\alpha_i \geq 0$

- ▶ **Convex combination:** $\mathbf{y} = \sum_{i=1}^N \alpha_i \mathbf{v}_i$ where $\sum_{i=1}^N \alpha_i = 1$ and $\alpha_i \geq 0$

Transform your data

- Let $\mathbf{x} \in \mathbb{R}^d$ be an input vector and $\mathbf{W} \in \mathbb{R}^{d \times m}$ be a transformation matrix, a simple linear transform is:

$$\mathbf{y} = \mathbf{W}^T \mathbf{x} \quad (1)$$

- Random Projections:** Project data to **low** dimensional space using Random Gaussian matrix

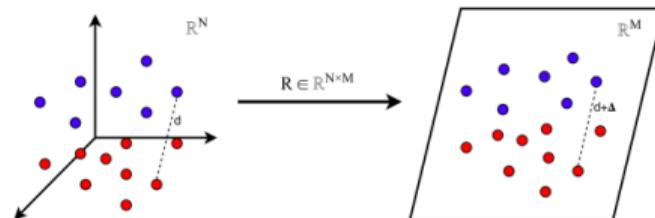


Figure: Random projections for dimensional reduction.

Types of Transforms



- ▶ Transforms can be categorised based on the **nature of transformation matrix**
 - ▶ **Analytical transform:** W is **analytically** determined
 - ▶ **Learned transform:** W is **learned** from the **data** itself
- ▶ Analytical transforms:
 - ▶ Data **independent** and universal
 - ▶ E.g. **Fourier** and **cosine** transforms where signal is decomposed onto sinusoidal waves of different frequencies
- ▶ Learned transforms:
 - ▶ Data **dependent**
 - ▶ W contains basis of a subspace or exemplars
 - ▶ E.g. **PCA** as **principle components** are learned from the data
 - ▶ **Dictionary learning:** Learn transformation matrices or dictionaries from data and decompose data onto the atoms of these dictionaries

Types of Transforms: Deep Learning



- Deep neural networks can be seen as **cascades** of learned transforms:

$$\mathbf{y} = \mathbf{W}_3 \mathbf{W}_2 \mathbf{W}_1 \mathbf{x} \quad (2)$$

But this is linear! (Deep dictionary learning¹)

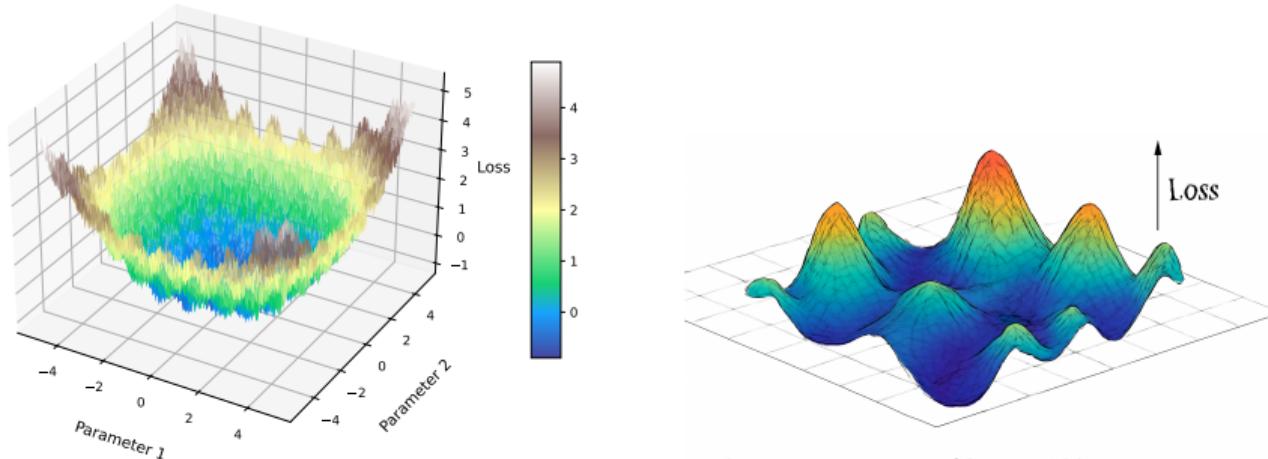
$$\mathbf{y} = \text{RELU}(\mathbf{W}_3 \text{RELU}(\mathbf{W}_2 \text{RELU}(\mathbf{W}_1 \mathbf{x}))) \quad (3)$$

- DNNs **learn** these transforms **simultaneously** using gradient descent
- Purpose is to learn embedding or **y**-space where classes are separable

¹ Thakur & Rajan, Deep archetypal analysis based intermediate matching kernel for bioacoustic classification, IEEE JSTSP 2019.

Loss Landscape

- ▶ A loss landscape represents how the loss (or error) of a neural network changes as we adjust its weights.



Loss Landscape of Neural Networks

Loss Landscape

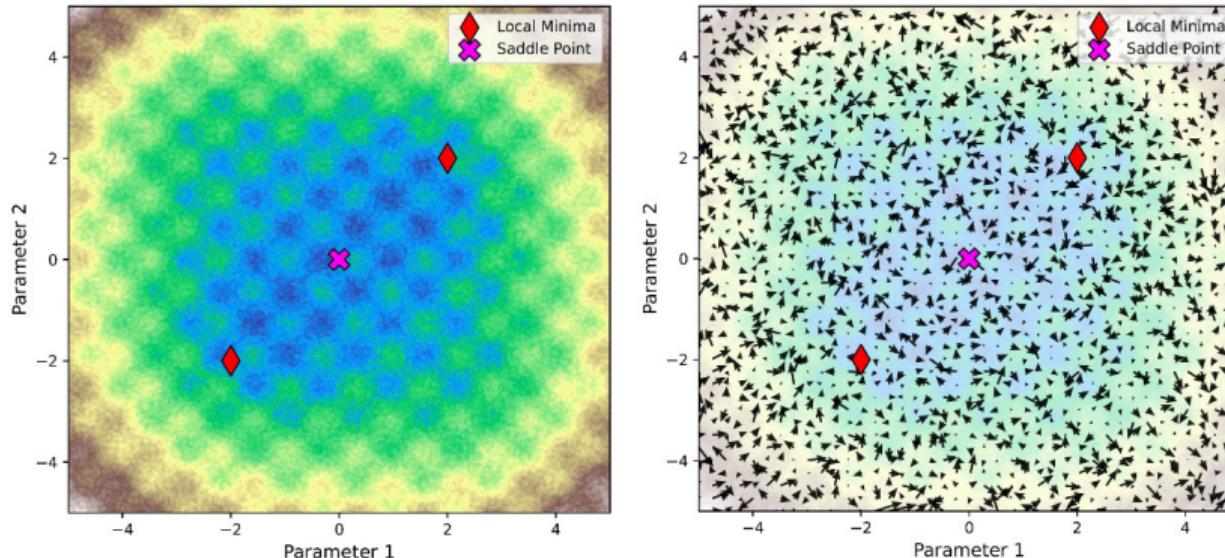
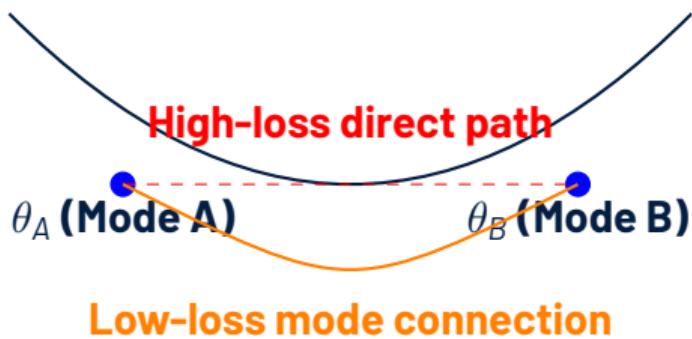


Figure: Relation between loss landscape and gradient field.

What is a Mode Connection?

- ▶ A **Mode Connection** is a **low-loss path** between two local optima (or **modes**) of a neural network.
- ▶ This path satisfies:
 - ▶ **Smooth Transition:** Every point along the path remains a viable local minimum.
 - ▶ **Loss Constraint:** The entire path maintains loss values **no worse** than the endpoints.



What is a Mode Connection?

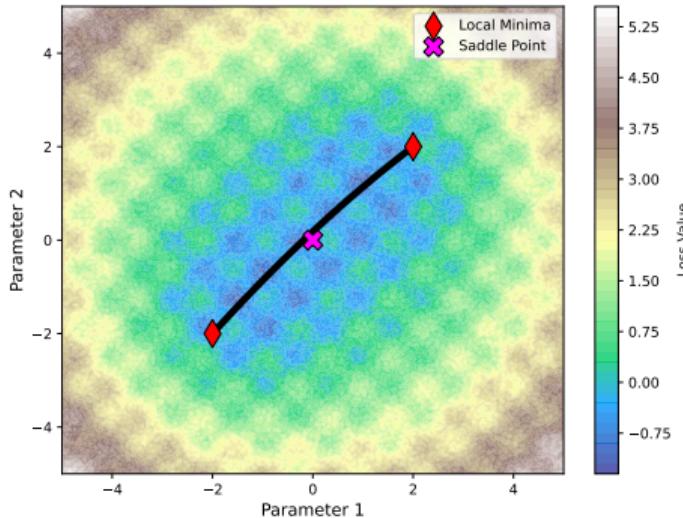


Figure: Mode connection in loss landscape.

Why is Mode Connectivity Important?



► **Understanding the Loss Landscape:**

- Challenges the assumption that local minima are isolated.
- Reveals the underlying structure of neural network loss surfaces.

► **Finding Flatter Minima for Better Generalization:**

- Improves **robustness** by making models less sensitive to **small input changes**.
- Enhances stability by reducing sensitivity to **weight perturbations**.
- Leads to **smoother decision boundaries**, improving test-time performance.

► **Providing Insights for Optimization:**

- Guides new training strategies for selecting better minima.
- Helps understand interactions in **joint optimization and model merging**.

- ▶ **Interpolation** between two optima using a direct linear path in parameter space.

$$\theta(t) = (1 - t)\theta_A + t\theta_B, \quad t \in [0, 1]$$

- ▶ Often passes through **loss barriers**
- ▶ Observed when loss landscape is **approximately convex** and models are trained similarly

Non-Linear Mode Connectivity

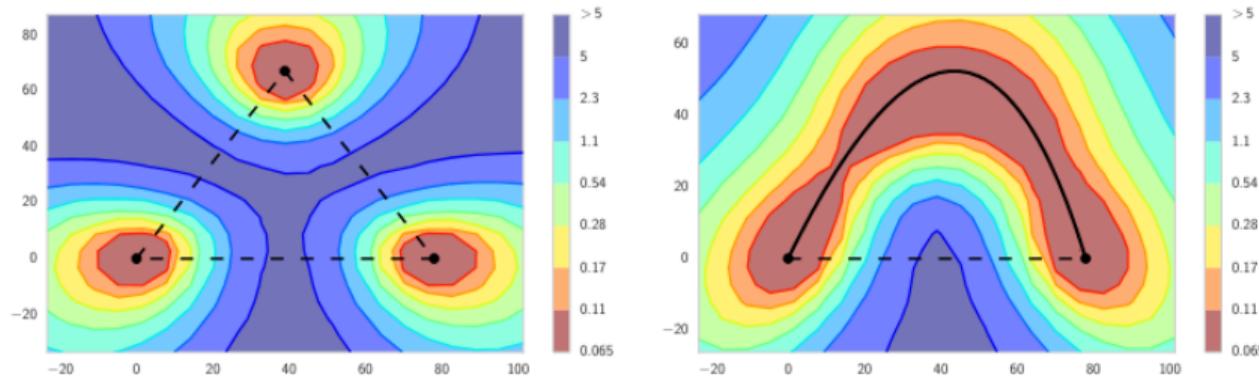


Figure: Loss Landscape of Resnet-164 trained on CIFAR-100.

Garipov et al., *Loss surfaces, mode connectivity, and fast ensembling of DNNs*, NeurIPS (2018).

- ▶ **Bezier Curves:** Trainable Parametric curves that can act as mode connections

$$\theta_\phi(t) = (1-t)^2 \theta_A + 2(1-t)t\phi + t^2 \theta_B, \quad t \in [0, 1]$$

- ▶ ϕ is the learnable **control point**, optimized to minimize loss along the curve.
- ▶ **Training the Bézier Curve**
 - ▶ Instead of direct weight interpolation, we **optimize** ϕ to ensure a **low-loss path**.
 - ▶ Training objective:

$$\mathcal{L}_{\text{Bézier}} = \mathbb{E}_t [\mathcal{L}(\theta_\phi(t))]$$

- ▶ We **sample values of** t , compute loss at those points, and optimize ϕ to minimize the loss along the curve.

Non-Linear Mode Connectivity in Action

- Respiratory deterioration prediction model

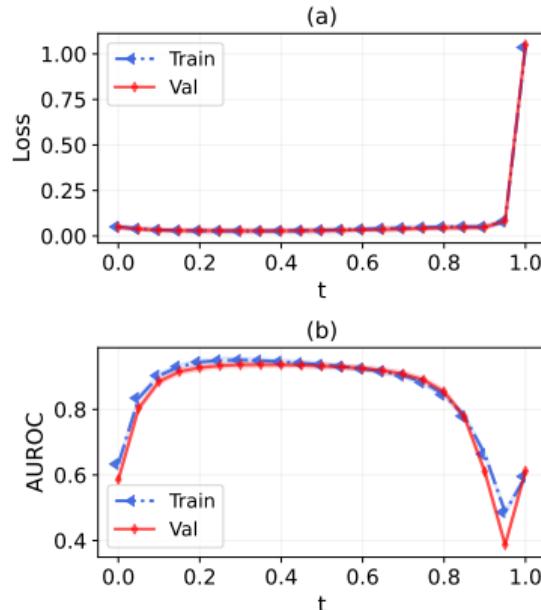


Figure: Mode connection between an optimum and a random point in parameter space.

Table of Contents

Preliminary Concepts

Transforms

Loss Landscape & Mode Connections

Multi-tasking

Adaptive Parameter Optimisation

Task Groupings

Gradient-based Meta-Learning

Federated Learning

Data Democratisation & Healthcare

Data Encoding and Condensation

Uncertainty in LLMs



Multi-tasking Architectures

- ▶ Multi-tasking architectures are characterised by:
 - ▶ A **shared** feature extractor “trunk”
 - ▶ Tasks-specific layers working on **same** feature representation

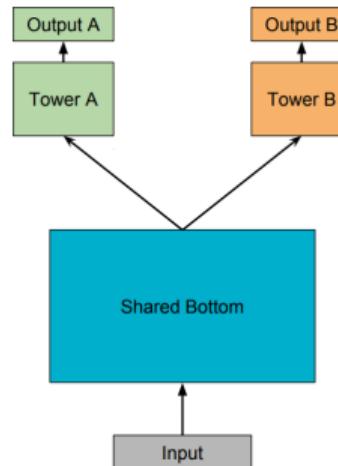


Figure: Illustration of a typical multi-tasking architecture.

Multi-tasking: Under the hood



- ▶ \mathbf{x} : input examples. (y_1, y_2, \dots, y_n) : corresponding task labels
- ▶ Model $f_{\theta, \phi_i}()$ outputs task-specific predictions: $\mathbf{p}_1, \mathbf{p}_2, \dots, \mathbf{p}_n$
- ▶ θ : shared parameters, ϕ_i : task-specific parameters for i th task
- ▶ Loss function is an **aggregate** of task-specific losses:

$$\mathcal{L} = \sum_{i=1}^n \alpha_i \mathcal{L}_i(y_i, \mathbf{p}_i), \tag{4}$$

where α_i is the coefficient for task i .

Multi-tasking: Under the hood

- ▶ Computing gradients for shared parameters:

$$\nabla_{\theta} \mathcal{L} = \sum_{i=1}^n \alpha_i \nabla_{\theta} \mathcal{L}_i(y_i, \mathbf{p}) \quad (5)$$

- ▶ α_i are usually constrained to be the convex coefficients

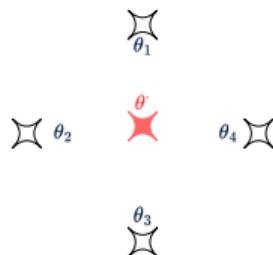
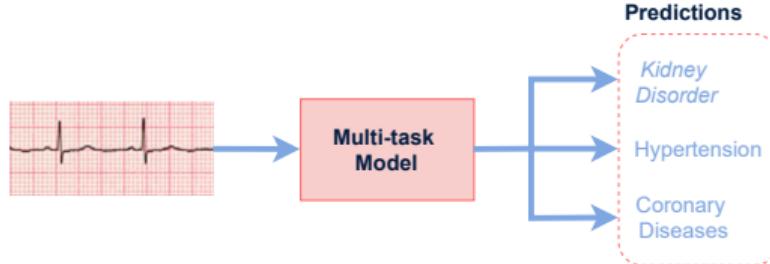


Figure: Deviation between optimal parameters and shared parameters learned by MTL.

Why do we need multi-tasking?



- ▶ Predicting **multiple outcomes** for a patient from a single input:
 - ▶ **Avoids** firing **multiple models**
 - ▶ **Lesser** storage complexity
- ▶ Regularises the training
 - ▶ **Avoids** over-fitting
 - ▶ Root cause of this regularisation is learning **common representation**
- ▶ Another way to look at it:
 - ▶ MTL adds noise to noisy gradients
 - ▶ Further increasing the chances of arriving at wider local minima

Challenges in Multi-tasking: Shared Representation



- ▶ Is it always possible to learn a **good** shared representation?
 - ▶ Of-course not!
- ▶ A **good** shared representation could be obtained if tasks are **similar**
 - ▶ Our and machine's perception of task similarity could be different!
- ▶ **Less similar** tasks result in shared parameters that are **not be helpful** to any task

Challenges in Multi-tasking: Optimisation Issues

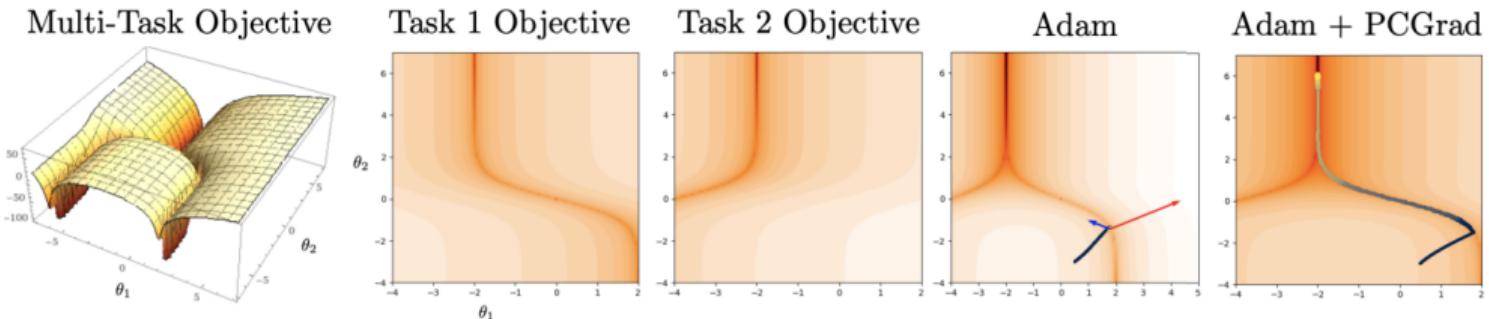


Figure: Gradient conflicts

Source: Yu et al., Gradient Surgery for Multi-Task Learning, Neurips 2020.

Challenges in Multi-tasking: Optimisation Issues

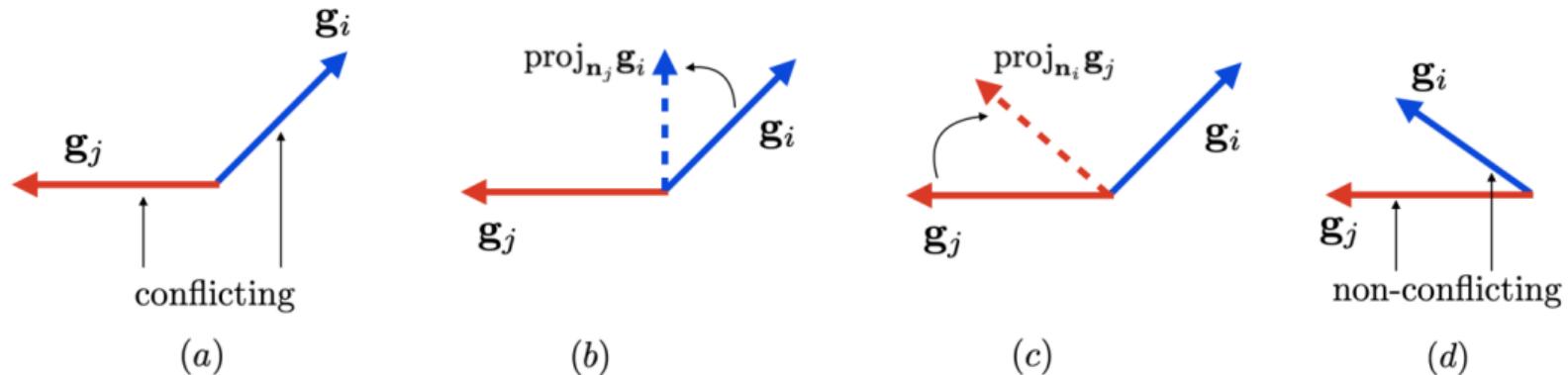


Figure: Gradient conflicts

Source: Yu et al., Gradient Surgery for Multi-Task Learning, Neurips 2020.

Challenges in Multi-tasking: Label Availability

- ▶ Availability of **all labels** for an example is **rare**
- ▶ In healthcare informatics, it's possible that a patient may not exhibit all the outcomes that we are modelling
- ▶ MTL can **only** use a **subset** of all the available data

Lazy Training²

- Minimal or no deviation in **majority** of parameters

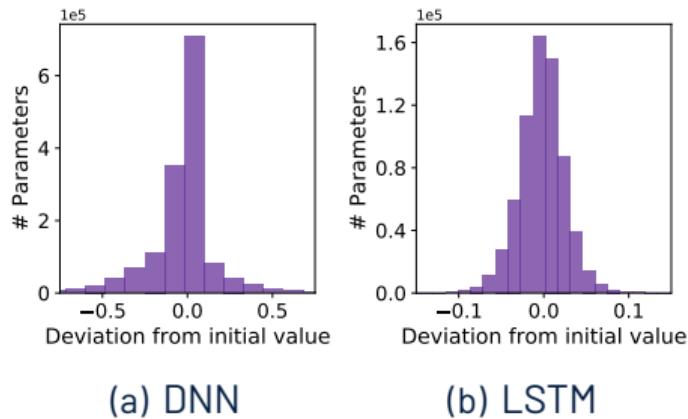


Figure: Histograms depicting deviation of trained parameters from their initial values in (a) 5-layered DNN trained on FashionMNIST dataset and (b) LSTM based model trained on MIMIC-III dataset for in-hospital mortality prediction.

² Chizat et al., On lazy training in differentiable programming, Neurips 2019.

Adaptive Parameter Optimisation: An Overview

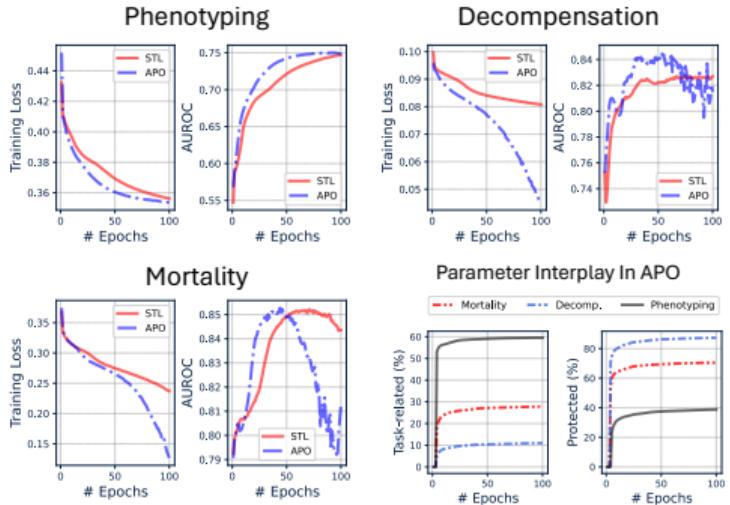
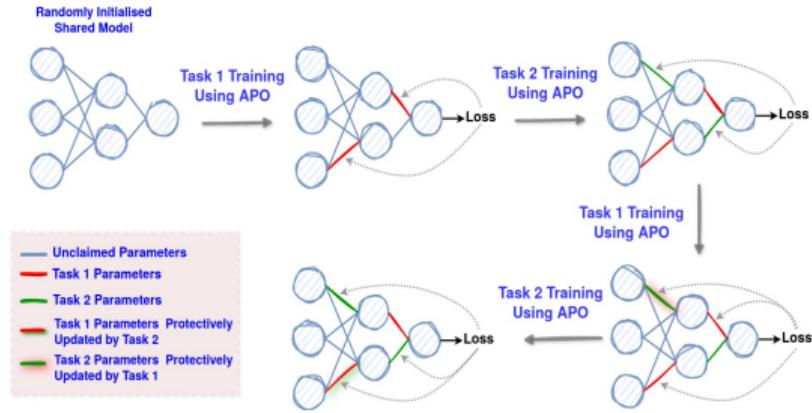


Figure: Illustration of APO³ in action.

³

Thakur et al., Information Transfer Across Clinical Tasks via Adaptive Parameter Optimisation, AISTATS 2025.

Task Groupings to Alleviate Negative Transfer

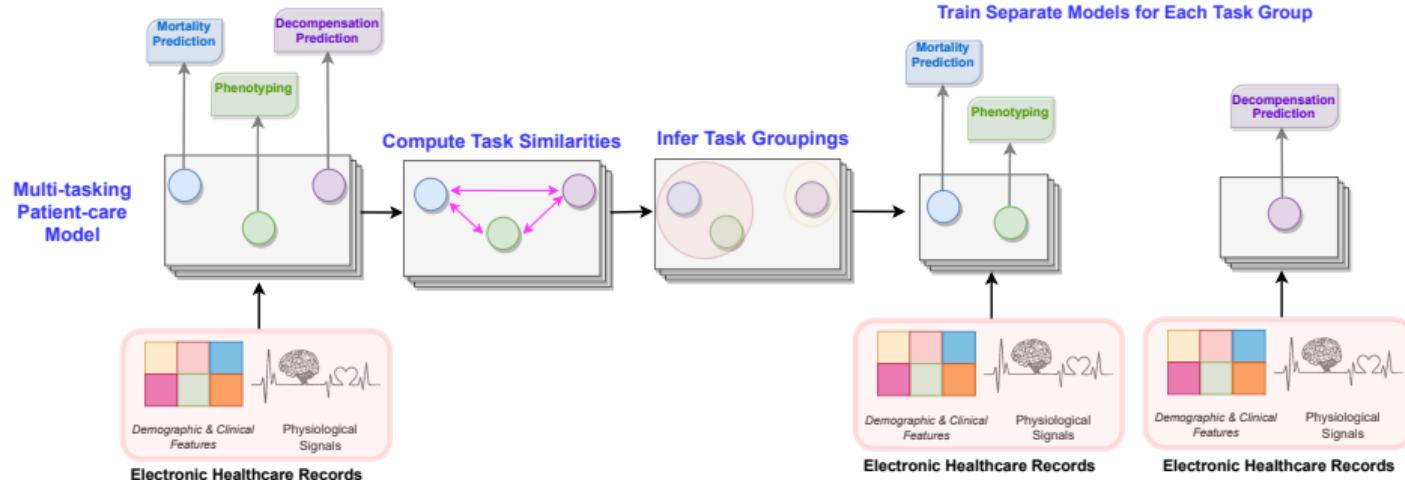


Figure: A typical task grouping framework⁴.

⁴

Thakur et al., Efficient Task Grouping Through Sample-wise Optimisation Landscape Analysis, IEEE TPAMI 2025.

Sample-wise optima



- ▶ $f_\theta(\cdot)$: neural network parameterised by θ , and θ_0 is the initial random state of the model
- ▶ i -th sample: $(\mathbf{x}_i, y_i^1, y_i^2, \dots, y_i^t)$ in dataset \mathcal{D}
- ▶ Task-specific sample-wise optima θ_i^{t*} :

$$\theta_i^{t*} = \theta_0 - \nabla_{\theta} \ell(f_\theta(\mathbf{x}_i), y_i^t)$$

- ▶ Sample-wise optima θ_i^* :

$$\theta_i^* = \frac{1}{T} \sum_{t=1}^T \theta_i^{t*}$$

Sample-wise Convergence For Task Groupings

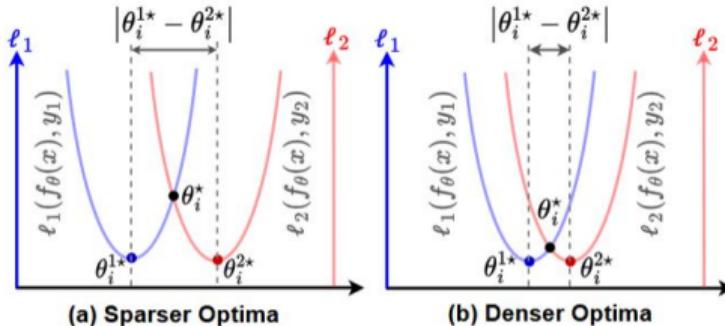


Figure: Illustration of the nature of shared global optima θ^* as a function of the density of task-specific sample-wise local optima, θ_1^* and θ_2^* for tasks 1 and 2, respectively. In comparison to sparser local optima **(a)**, denser local optima **(b)** results in θ^* that leads to better average loss across both tasks, $\mathcal{J} = (\ell_1 + \ell_2)/2$, as desired in multi-tasking and other information sharing frameworks.

Sample-wise Convergence For MTL

- Average loss bound:

$$\mathcal{J} \leq \frac{T^3}{2n} \sum_{i=1}^n \psi_{\theta_i}^2 + \frac{n^3}{2} \Psi_{\theta^*}^2,$$

$$\psi_{\theta_i} = \frac{\sqrt{\mathcal{H}}}{T^2} \sum_{j=1}^T \sum_{k=1}^T \|\theta_i^{j*} - \theta_i^{k*}\|_1,$$

$$\Psi_{\theta^*} = \frac{\sqrt{\mathcal{H}}}{n^2} \sum_{j=1}^n \sum_{k=1}^n \|\theta_j^* - \theta_k^*\|_1.$$

- SCA based task affinity between task τ_1 and τ_2 :

$$a_{\tau_i, \tau_j} = \frac{1}{n} \sum_{i=1}^n \|\theta_i^{\tau_i*} - \theta_i^{\tau_j*}\|_1,$$

Results: Celeb-A dataset



Table: Performance of different task grouping methods on Celeb-A dataset in terms of (a) total absolute error and (b) average training time (RTX A6000 GPU hours).

(a) Total Absolute Error

Method ↓	Splits			
	2	3	4	5
Random	50.1 ± 0.085	49.97 ± 0.13	49.75 ± 0.08	49.73 ± 0.11
TAG	49.66 ± 0.095	49.55 ± 0.08	49.48 ± 0.066	49.46 ± 0.07
HOA	49.76 ± 0.08	49.74 ± 0.43	49.72 ± 0.1	49.69 ± 0.09
SCA(Prop.)	49.6 ± 0.24	49.58 ± 0.07	49.48 ± 0.05	49.44 ± 0.09

(b) Average Training Time

Method ↓	Splits			
	2	3	4	5
Random	3.52 ± 0.02	5.24 ± 0.02	6.98 ± 0.01	8.62 ± 0.02
TAG	5.4 ± 0.03	7.17 ± 0.04	8.61 ± 0.03	10.45 ± 0.02
HOA	57.6 ± 0.08	59.5 ± 0.06	61.2 ± 0.04	62.7 ± 0.06
SCA(Prop.)	3.61 ± 0.02	5.31 ± 0.03	7.06 ± 0.02	8.71 ± 0.01

Results: Celeb-A dataset

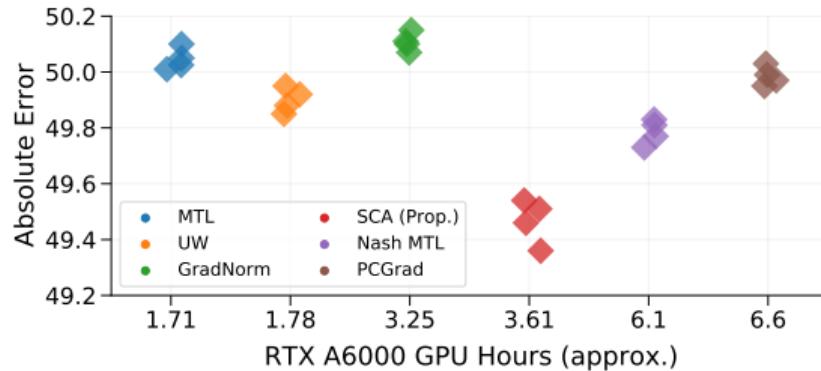


Figure: Performance and computational efficiency of the proposed SCA-based framework in comparison to prominent MTL methods on Celeb-A dataset.

Table of Contents

Preliminary Concepts

Transforms

Loss Landscape & Mode Connections

Multi-tasking

Adaptive Parameter Optimisation

Task Groupings

Gradient-based Meta-Learning

Federated Learning

Data Democratisation & Healthcare

Data Encoding and Condensation

Uncertainty in LLMs



- ▶ **Interpretation:** Training a model using gradients computed by other models
- ▶ **Purpose:** Learning a common model across multiple tasks
 - ▶ Global/common model can be quickly adapted to new unseen tasks
 - ▶ Proposed to mimic the learning in humans
- ▶ **Why are we interested in meta-learning?:**
 - ▶ Provides common/shared model across multiple tasks
 - ▶ Alleviate **example-labels requirement** of MTL
 - ▶ Information sharing across tasks in an example-independent manner

First Order Meta-learning: REPTILE⁵

- ▶ $f_\theta()$: Global model parameterised by θ
- ▶ $\mathcal{D}_t = \{\mathbf{x}_i, y_i\}_{i=1}^n$: Task t dataset

for $t \leftarrow 1 : T$ **do**

$$\theta_t = \theta$$

$\mathcal{B} \leftarrow \text{SAMPLE-BATCHES}(\mathcal{D}_t)$

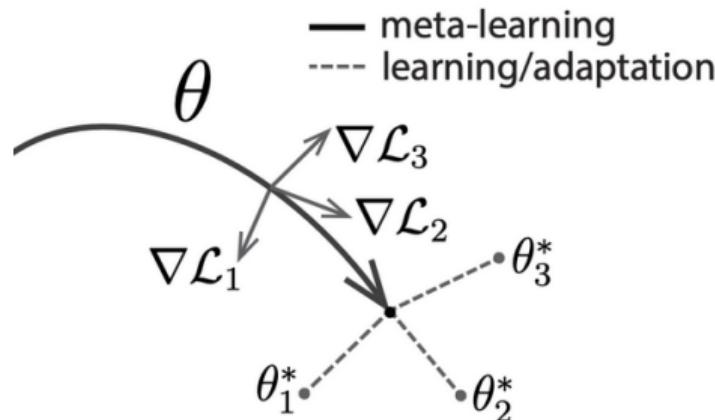
for any $(\mathbf{b}, \mathbf{l}) \in \mathcal{B}$ **do**

$$L = \mathcal{L}(f_{\theta_t}(\mathbf{b}, \mathbf{l}))$$

$$\theta_t = \theta_t - \alpha \nabla_{\theta_t} L$$

$$\phi = \sum_{t=1}^T (\theta_t - \theta)$$

$$\theta = \theta + \alpha \phi$$



Multi-tasking to Meta-learning



- ▶ All inner-loop updates can be summaries as one single gradient update
- ▶ Meta-grad for task t , $\phi_t = -(\theta - \theta_t) = \nabla_{\theta_t} L$
- ▶ REPTILE update can be written as:

$$\nabla \mathcal{L}_\theta = \sum_{t=1}^T \nabla_{\theta_t} L = \sum_{t=1}^T \alpha_i \nabla_{\theta_t} L \quad (6)$$

- ▶ REPTILE/MAML, MTL and **FedAvg** are identical w.r.t. optimisation

Table of Contents

Preliminary Concepts

Transforms

Loss Landscape & Mode Connections

Multi-tasking

Adaptive Parameter Optimisation

Task Groupings

Gradient-based Meta-Learning

Federated Learning

Data Democratisation & Healthcare

Data Encoding and Condensation

Uncertainty in LLMs

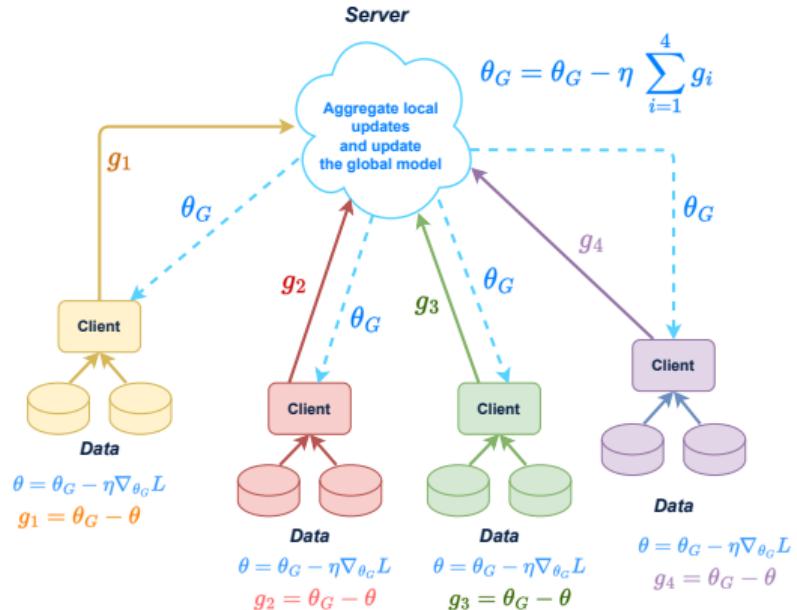
Data Privacy Constraints in Healthcare



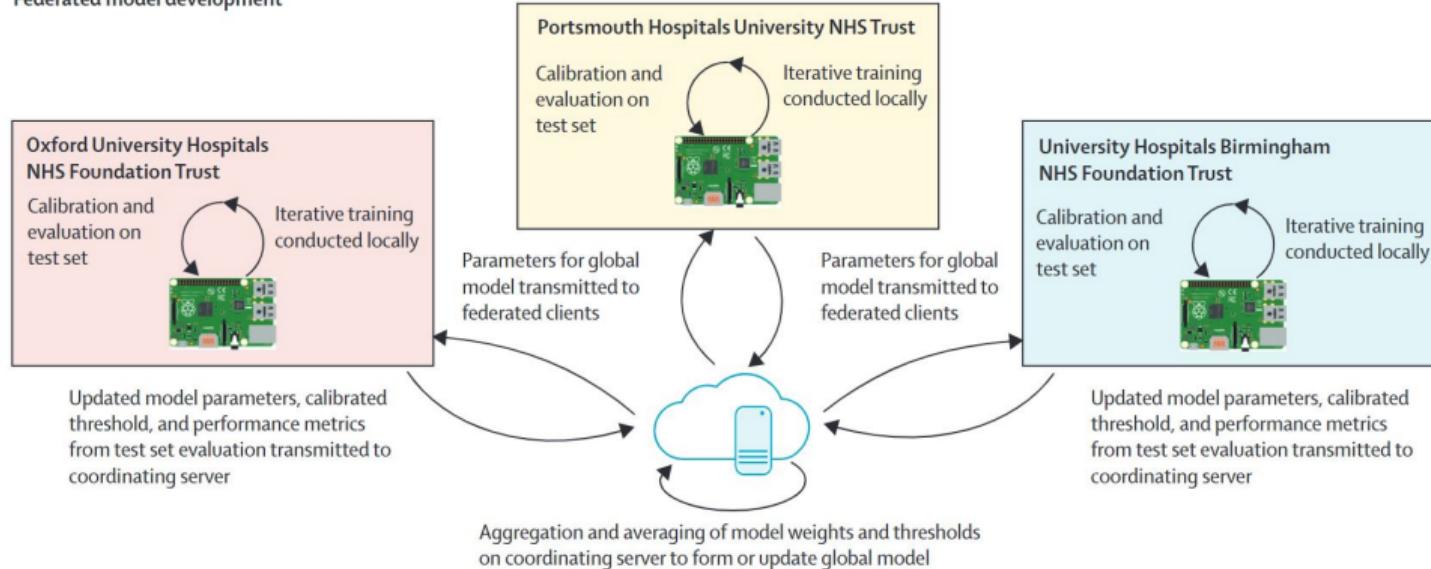
- ▶ Healthcare data is **sensitive**, contains patients' **private** information
- ▶ Protected by legislation such as **GDPR** and **DATA PRIVACY ACTS**
- ▶ Data is **distributed** across multiple sites such as the NHS Trusts, cannot be moved across sites
- ▶ **Global** trends may not be captured without the analysis of complete data
- ▶ Generalised predictive models should be a **function** of datasets distributed across sites

Federated Learning

- **Distributed** training of deep models using **FedSGD** or **FedAvg**
- **SGD**: $\theta' = \theta - \eta \nabla_{\theta} \mathcal{L}(f_{\theta}(\mathbf{x}), \mathbf{y})$
 - (\mathbf{x}, \mathbf{y}) : Data, $f_{\theta}()$: Neural network parameterised by θ

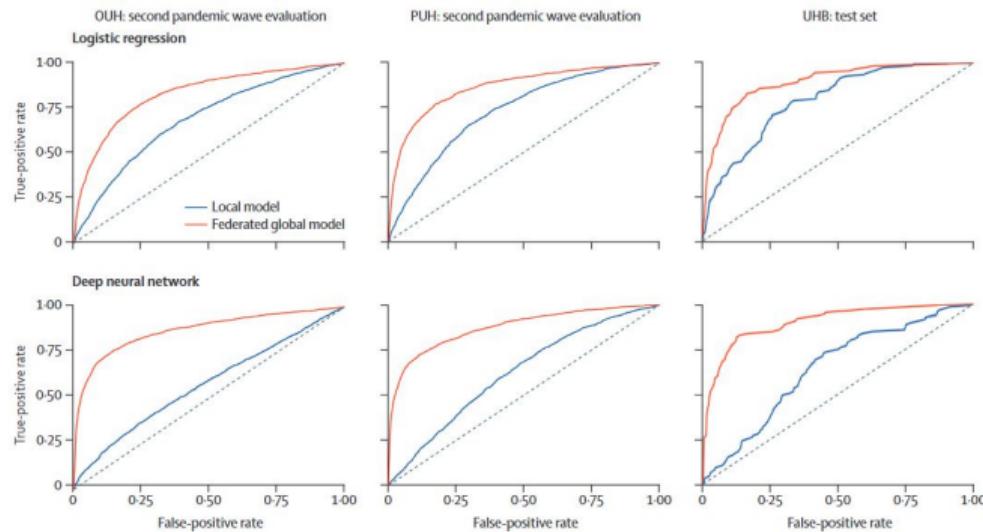


Federated model development



⁶

Soltan et al., A scalable federated learning solution for secondary care using low-cost microcomputing: privacy-preserving development and evaluation of a COVID-19 screening test in UK hospitals, *The Lancet Digital Health* 2023.



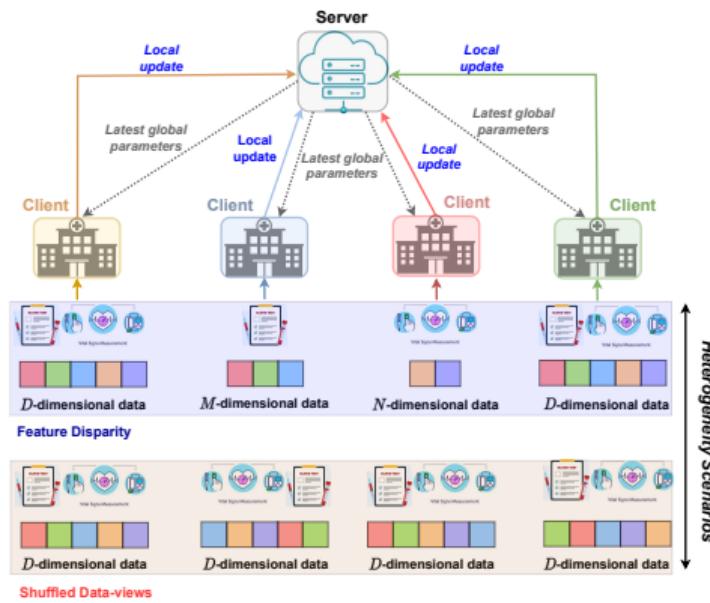
⁷

Soltan et al., A scalable federated learning solution for secondary care using low-cost microcomputing: privacy-preserving development and evaluation of a COVID-19 screening test in UK hospitals, *The Lancet Digital Health* 2023.

Federated Learning: A Perfect Solution?



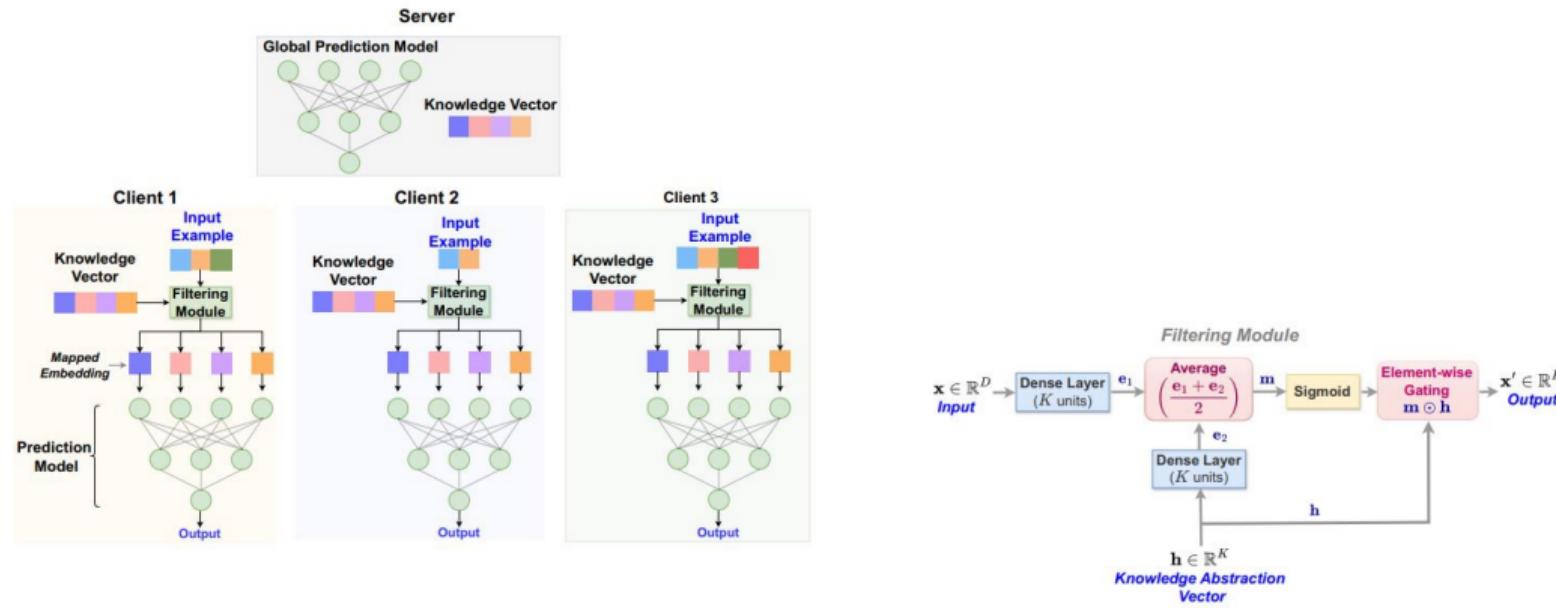
- ▶ **Data-view heterogeneity:** Massive pre-processing is required to align data-views across clients



Federated Learning under Data-view Heterogeneity



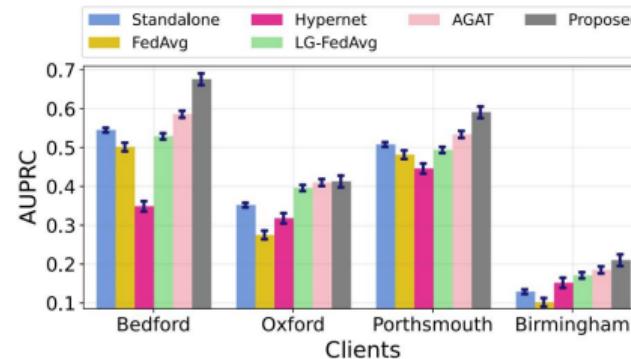
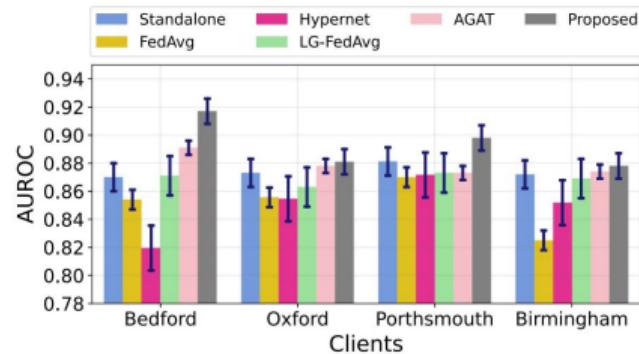
Federated Learning under Data-view Heterogeneity⁸



⁸

Thakur et al., Knowledge abstraction and filtering based federated learning over heterogeneous data views in healthcare, npj Digital Medicine 2024.

Federated Learning under Data-view Heterogeneity



Federated Learning under Data Heterogeneity

- ▶ FedAvg or FedSGD may struggle under **non-IID** settings
- ▶ **Client heterogeneity** may force local models to **deviate far away** from each other

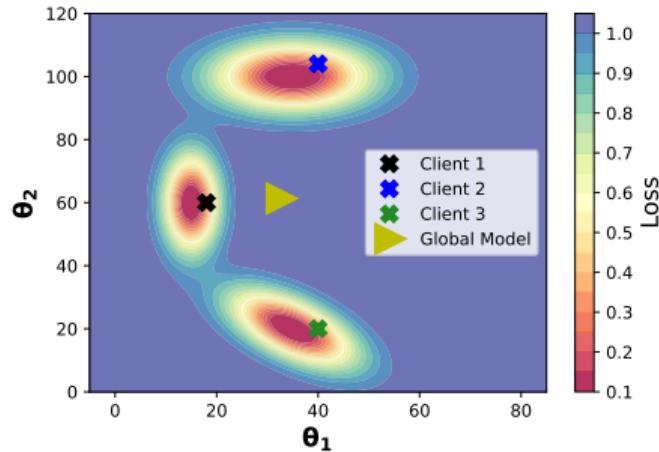
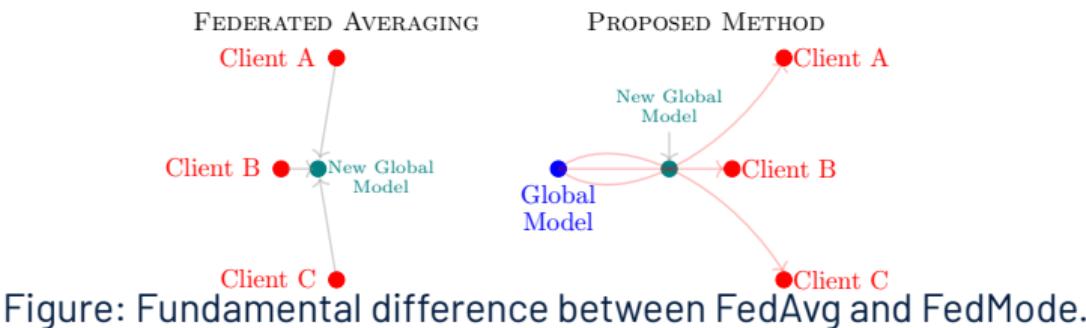


Figure: Global model drifts to high loss regions under client heterogeneity.

Optimising FL with Mode Connectivity

Key Idea: Improving Global Model Aggregation⁹

- ▶ **Aligning Client Models:** Leverage mode connections between the global model and clients' local models
- ▶ **Global Model Selection:** Choose the new global model **near the intersection** of these mode connections
- ▶ **Improved Fairness:** Ensures the global model resides in **a shared low-loss region** that benefits all clients



⁹

Thakur et al., Optimising Clinical Federated Learning through Mode Connectivity-based Model Aggregation, AISTATS 2025.

Optimising FL with Mode Connectivity

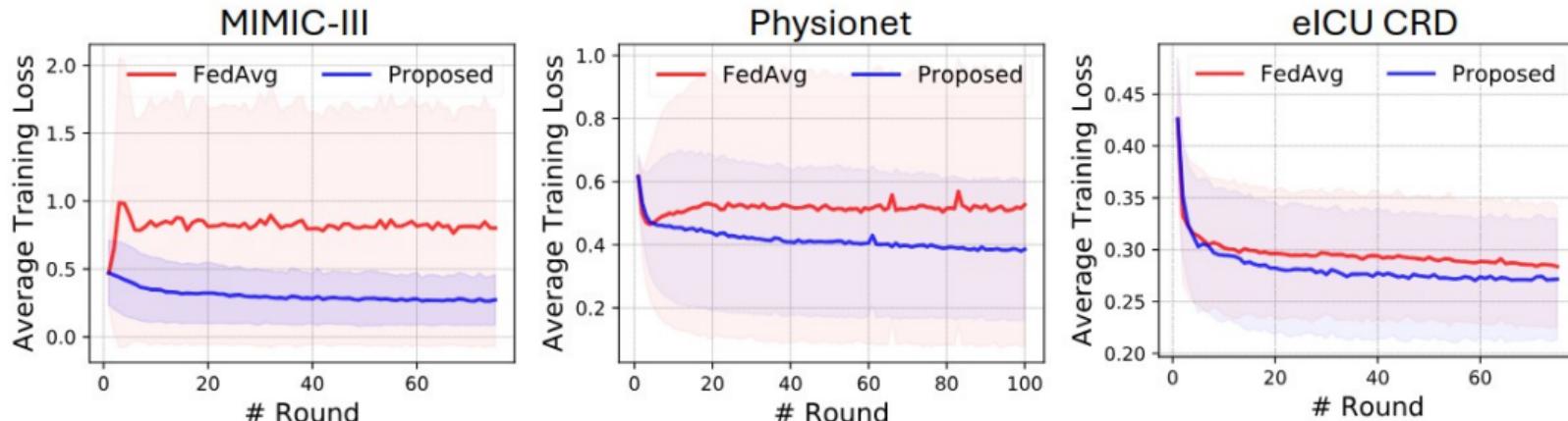


Figure: Training dynamics on 15 clients (MIMIC-III, Physionet) and 100 clients (eICU-CRD). MIMIC-III and Physionet exhibit higher heterogeneity.

Optimising FL with Mode Connectivity: Client Fairness

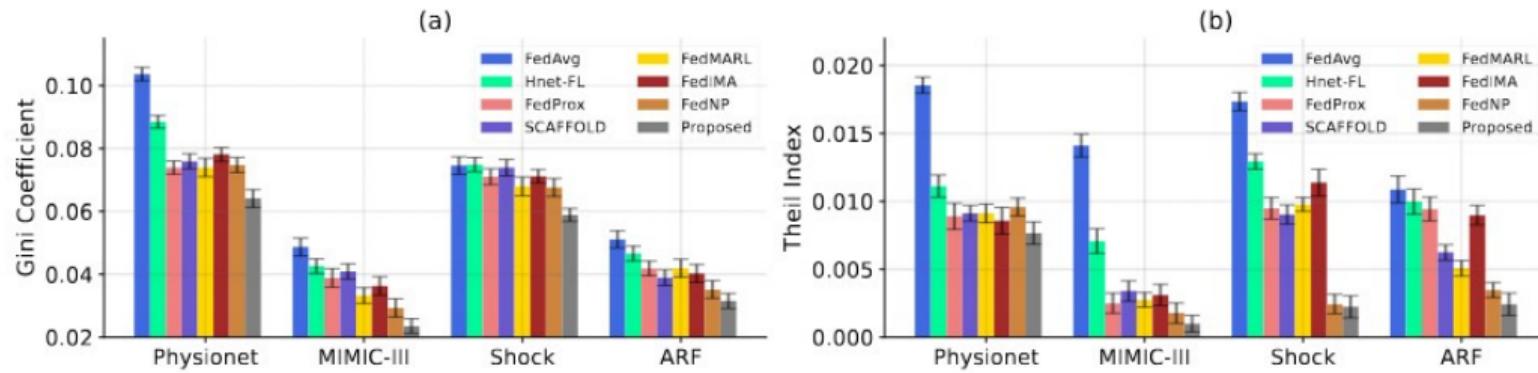


Figure: Client fairness across different FL methods.

Federated Learning: A Perfect Solution?

- ▶ **Gradient Inversion:**¹⁰ Original samples can be reconstructed from the gradient updates
- ▶ Some success but can have large implications



¹⁰ Huang, Yangsibo, et al. "Evaluating gradient inversion attacks and defences in federated learning." NeurIPS, 2021.

Federated Learning: A Perfect Solution?



- ▶ No true **data democratisation**
- ▶ All clients are forced to work towards **common goal**
- ▶ Can be difficult to convince institutions to participate

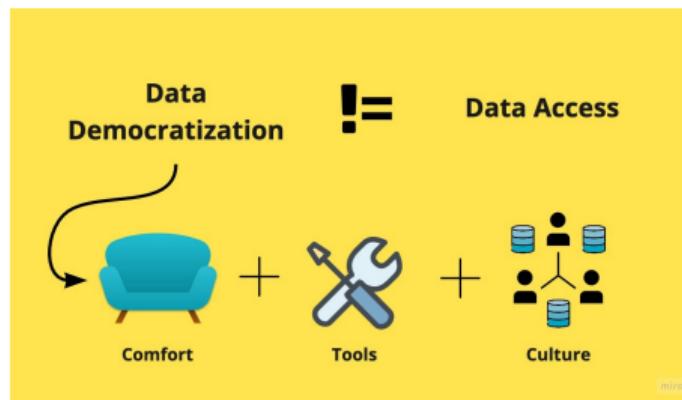


Table of Contents

Preliminary Concepts

Transforms

Loss Landscape & Mode Connections

Multi-tasking

Adaptive Parameter Optimisation

Task Groupings

Gradient-based Meta-Learning

Federated Learning

Data Democratisation & Healthcare

Data Encoding and Condensation

Uncertainty in LLMs



Data Democratisation in Clinical Machine Learning



- ▶ **Data Democratisation:** improving data access for **medical researchers**
- ▶ In the context of **Clinical ML:** data access for AI/ML researchers
- ▶ Most ML tools for healthcare informatics are neither robust nor validated
- ▶ **Data privacy laws** understandably inhibit data sharing
- ▶ Healthcare data democratisation could stimulate
 - ▶ **novel** clinical ML algorithms
 - ▶ **well-validated** models
 - ▶ **generic** models that are trained on multiple populations

Data Democratisation in Clinical Machine Learning



- ▶ Data Democratisation is possible **without violating** privacy constraints
- ▶ Key is to find a way to train models without **physical** data access
 - ▶ Extracting **imperceptible** semantic information from sensitive data
 - ▶ **Securely** transferring this information to the intended researchers
 - ▶ Researchers then use this information to train models

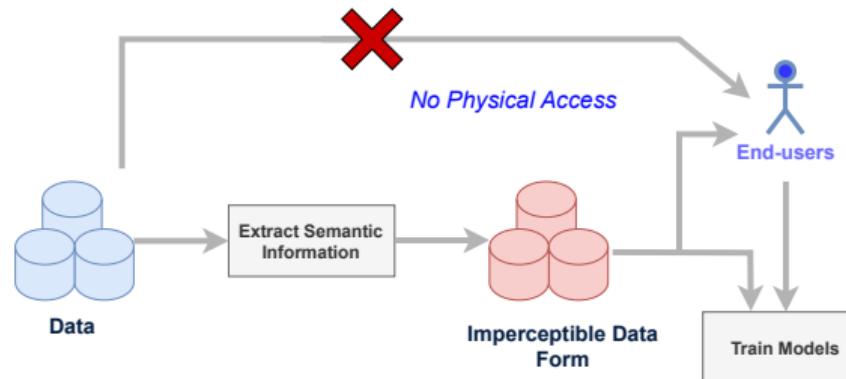


Table of Contents

Preliminary Concepts

Transforms

Loss Landscape & Mode Connections

Multi-tasking

Adaptive Parameter Optimisation

Task Groupings

Gradient-based Meta-Learning

Federated Learning

Data Democratisation & Healthcare

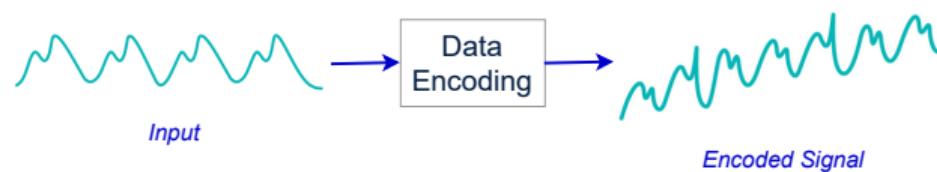
Data Encoding and Condensation

Uncertainty in LLMs



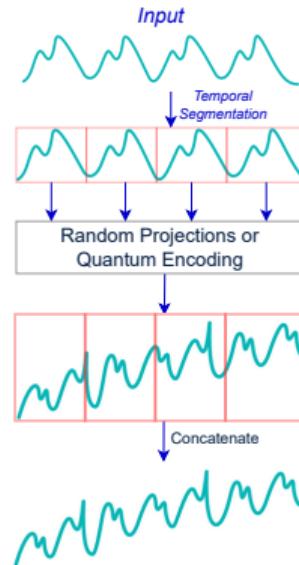
Data Encoding

- ▶ **FL:** Very restrictive, clients don't have **independence**
- ▶ **Data Encoding:** Transformation to make data "imperceptible" while preserving **semantics**
- ▶ Imperceptibility preserves data privacy
- ▶ Semantic preservation allows subsequent pattern analysis



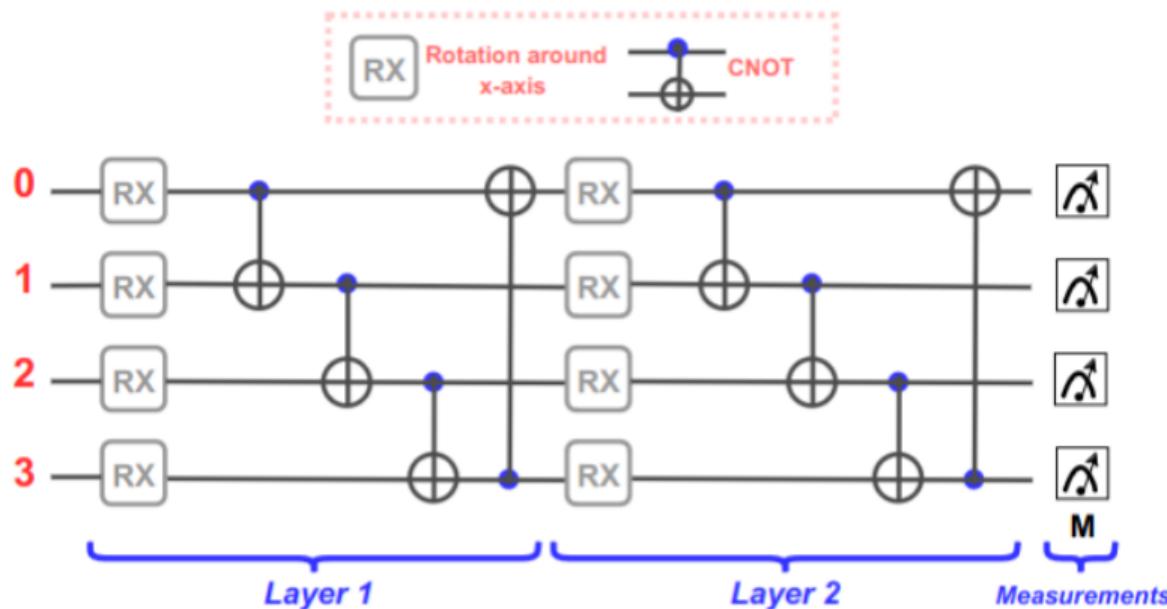
Data Encoding for Time-series Data

- ▶ Encoding framework:
 - ▶ Temporal segmentation of input signal
 - ▶ Apply data transformation on each segment
 - ▶ Concatenate the resulting segments
- ▶ Random Quantum Encoding and Random Projection as required transformations
- ▶ Each sample is a hash of its neighbours
- ▶ Loosing local structure to preserve global semantics



Thakur et al., *Data Encoding For Healthcare Data Democratisation and Information Leakage Prevention*, 2023.

Data Encoding for Time-series Data



Encoded Signals

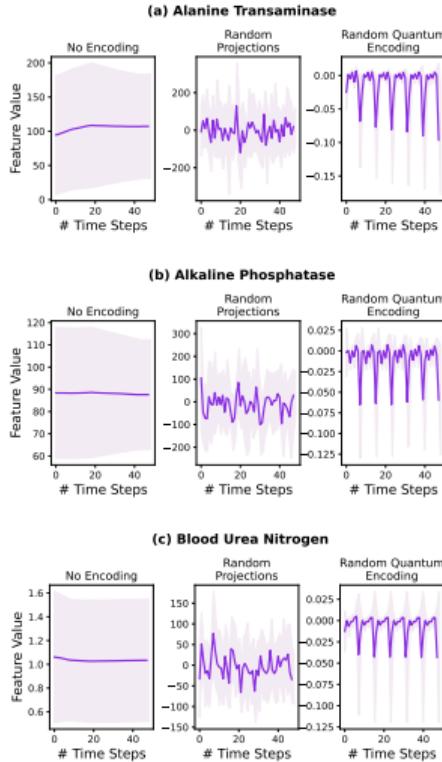
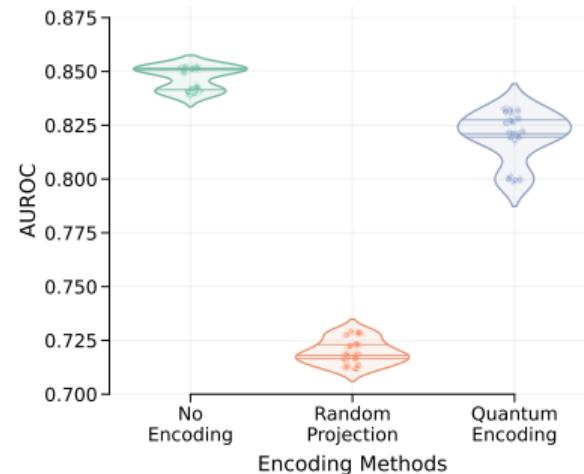
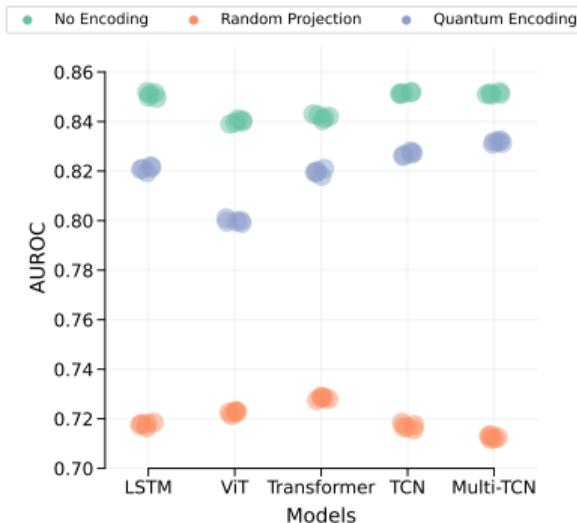


Figure: Average of 50 time series representing patients in the PhysioNet dataset.

Predictive Modelling on Encoded Data

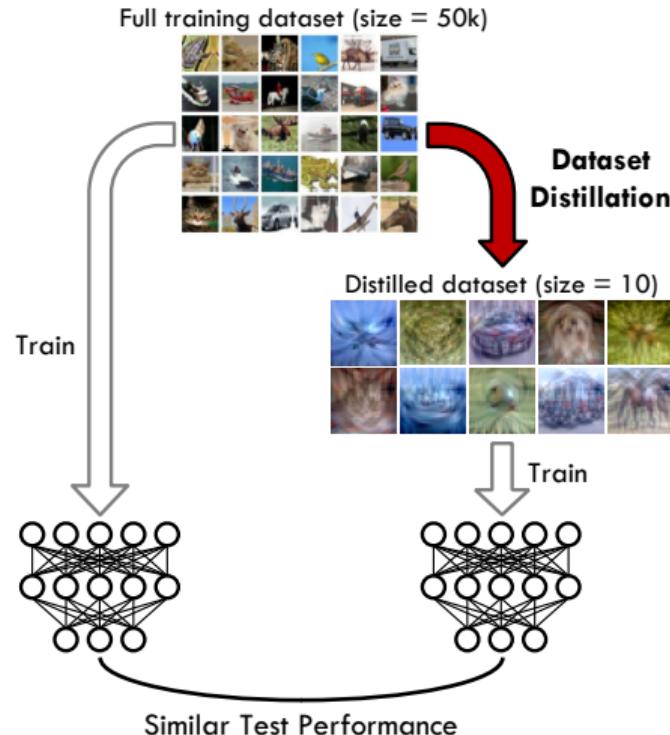
- **MIMIC-III** dataset for mortality prediction: 21,156 ICU stays
- **Models:** LSTM, Transformers, Vision Transformers (ViT), TCN & Multi-branch TCN



Mortality prediction on encoded MIMIC-III data.

Practicality outside academic circles?

Dataset Condensation

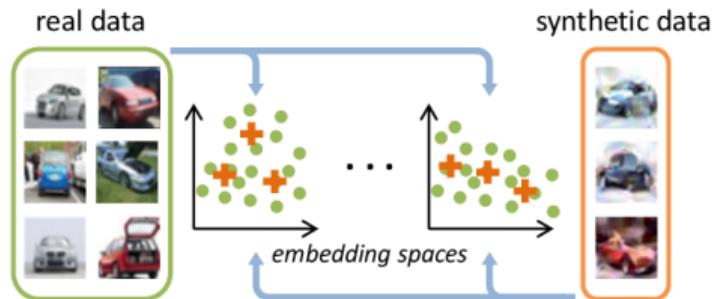


- ▶ **Privacy-preserving Nature:** Real Samples \rightleftharpoons Condensed Samples
 - ▶ No direct **one-to-one alignment**
 - ▶ Obfuscation via aggregation
 - ▶ Reduced **memorisation risks**
- ▶ **PRIVACY IS NO FREE LUNCH**
- ▶ **Differential Privacy** provides theoretical privacy guarantees

Dataset Condensation using Distribution Matching



- ▶ Aligning outputs of a function $f()$ generated for real \mathcal{D} and synthetic \mathcal{S} data
- ▶ $\mathbf{e}_1 = \frac{1}{N} \sum_{\forall \mathbf{x} \in \mathcal{D}} \mathbf{x}, \mathbf{e}_2 = \frac{1}{N} \sum_{\forall \mathbf{x} \in \mathcal{S}} \mathbf{x}$
- ▶ $\mathcal{L} = \|\mathbf{e}_1 - \mathbf{e}_2\|_2^2$
- ▶ $\mathbf{S} = \mathbf{S} - \eta \nabla_{\mathcal{S}} \mathcal{L}(\mathbf{e}_1, \mathbf{e}_2)$



Dataset Condensation: Time-series Data

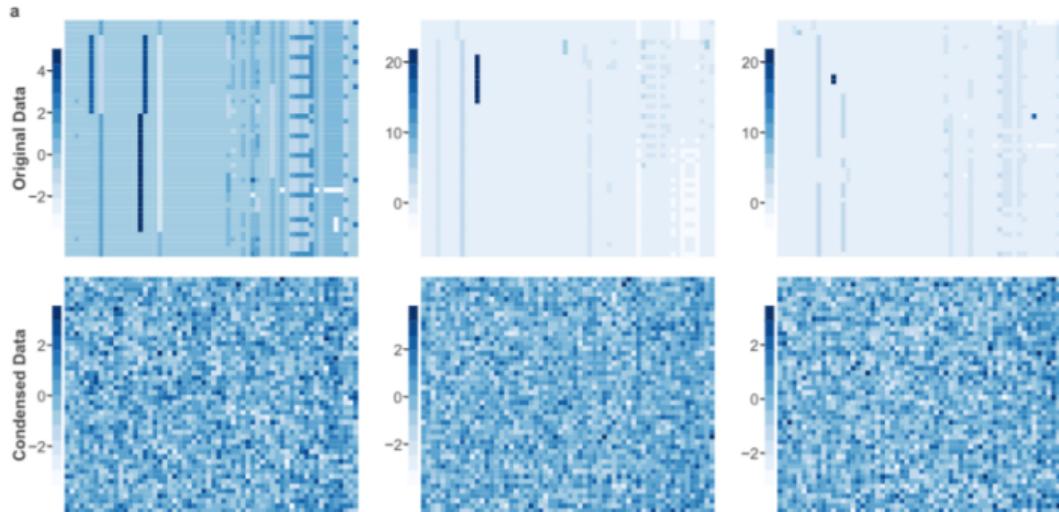
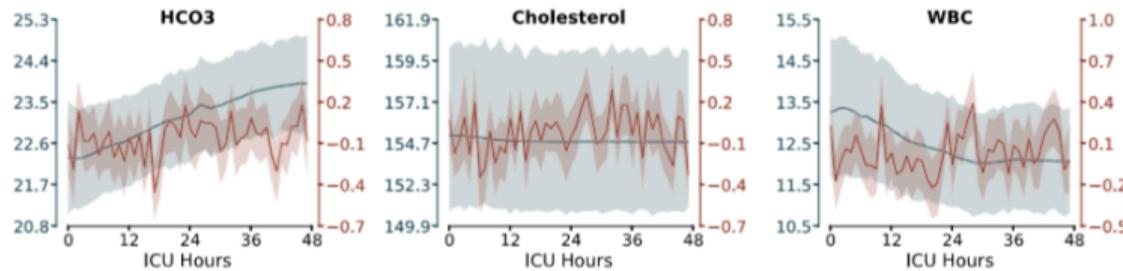


Figure: Heatmaps exhibiting changes in magnitude and trends of real and condensed time-series examples from MIMIC-III dataset.

Dataset Condensation: Time-series Data



Dataset Condensation: Time-series Data

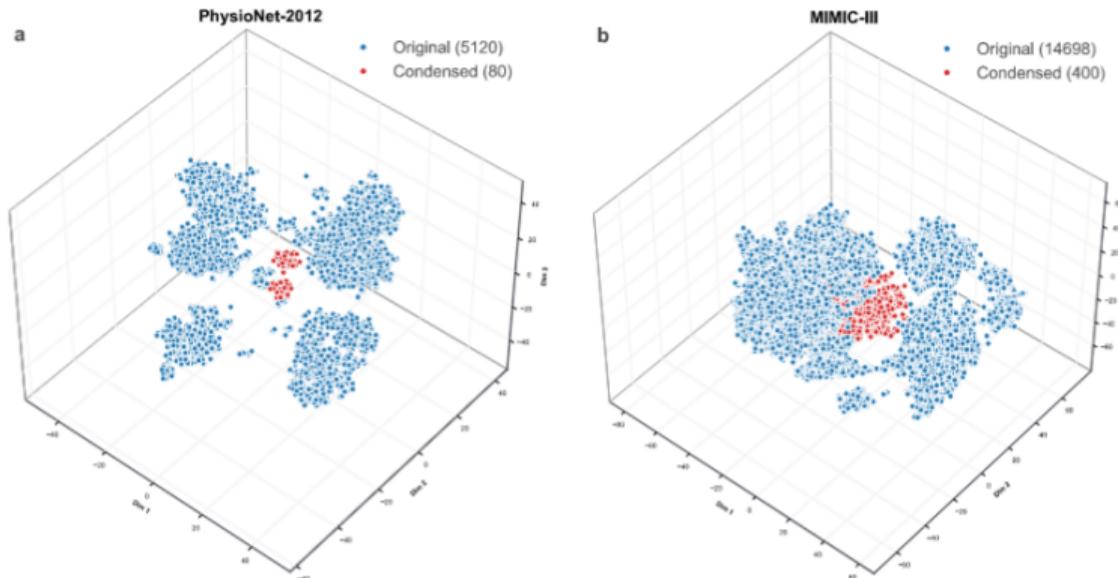


Figure: t-SNE representation of the average time-series examples.

Dataset Condensation: Predictive Performance

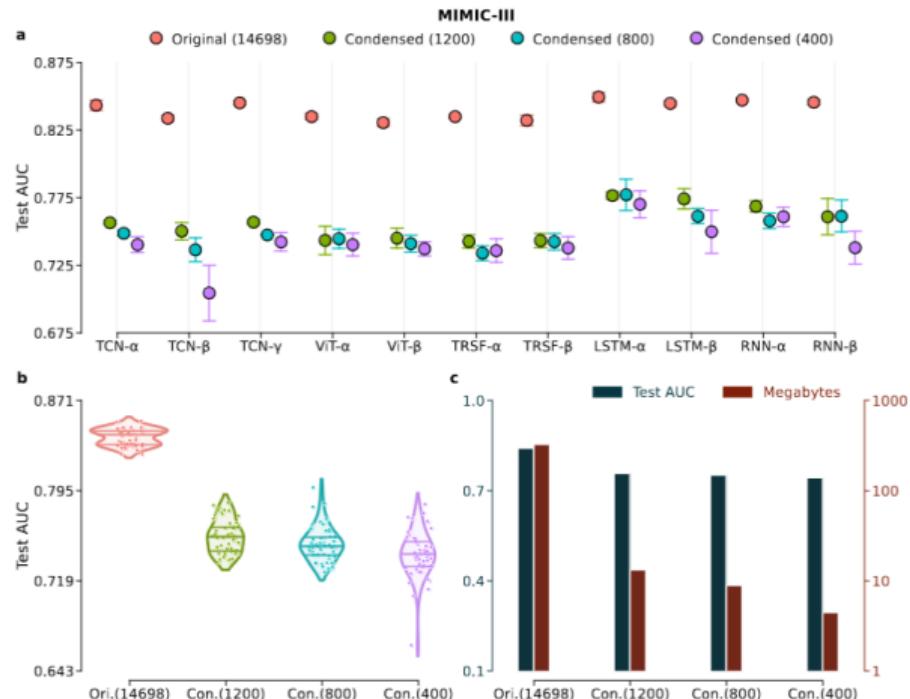


Figure: Performance on real and condensed data.

What is Trajectory Matching?

- ▶ State-of-the-art condensation methods that align synthetic data with real training trajectories
- ▶ Provide comparable performance to the real training data

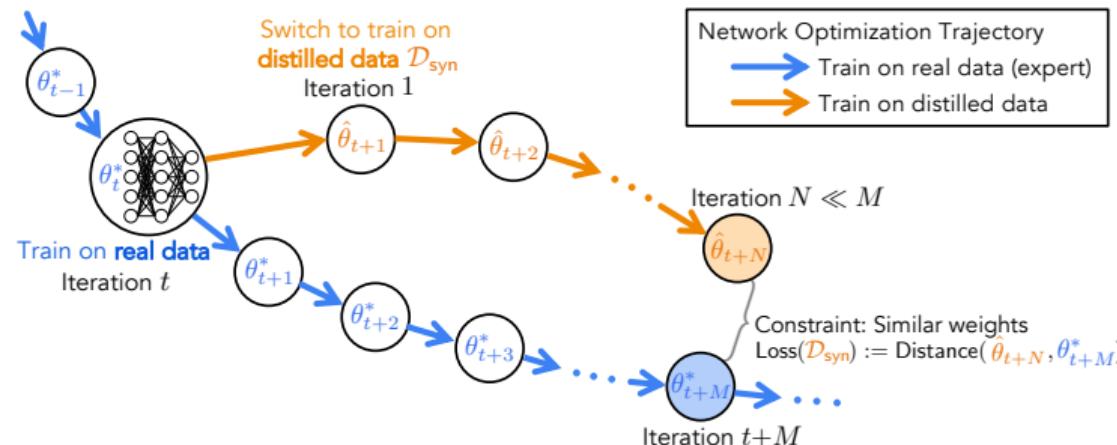
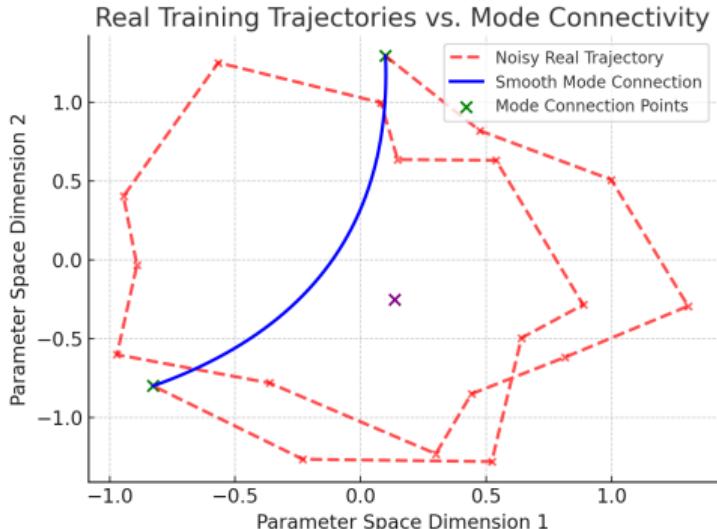


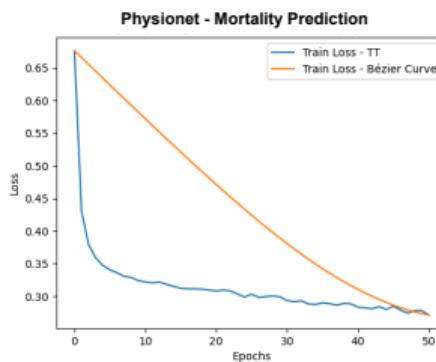
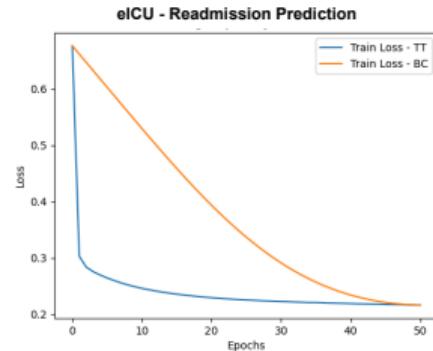
Figure: Illustration of trajectory matching strategy

Trajectory Matching: Drawbacks



- ▶ **HIGH MEMORY FOOTPRINT:** Requires storing full training trajectories across multiple models
- ▶ **NOISY AND UNSTABLE PATHS:** SGD trajectories are highly stochastic and non-smooth, reducing reliability for matching
- ▶ **AMBIGUOUS TEMPORAL STRUCTURE:** No clear notion of "early" or "late" stages in training, limiting interpretability and supervision quality

Mode Connections as Trajectory Surrogates¹¹



Theorem

Let $\mathcal{L} : \Theta \rightarrow \mathbb{R}$ be a β -smooth, lower-bounded loss function, $\theta_0 \in \Theta$ be a random initialization with loss $\ell_0 = \mathcal{L}(\theta_0)$, and θ_T be an SGD endpoint after K steps such that $\|\nabla \mathcal{L}(\theta_T)\| \leq \varepsilon$, where $\mathcal{L}(\theta_T) = \ell_T \ll \ell_0$. Assume $\gamma(t)$ denote the piecewise-linear interpolation of the SGD iterates $\{\theta_k\}_{k=0}^K$. The optimized Bézier curve and its curvature as $\theta^*(t) := \theta(t; \phi^*)$, $\kappa := 2\|\theta_0 - 2\phi^* + \theta_T\|$. Then,

(i) Average loss along the Bézier path is near-optimal:

$$\int_0^1 \mathcal{L}(\theta^*(t)) dt \leq \int_0^1 \mathcal{L}(\gamma(t)) dt + \frac{\beta \kappa^2}{240}. \quad (7)$$

(ii) Bézier path has lower and noise-free curvature compared to SGD trajectory:

$$\sup_{t \in [0, 1]} \|\theta^{**}(t)\| = \kappa, \quad \mathbb{E} \left[\sup_t \|\gamma''(t)\| \right] \geq \kappa + c \cdot \sigma_{\text{sgd}} \quad (8)$$

(iii) Model predictions along Bézier path remain close to those along SGD:

$$\sup_{x \in \mathcal{X}, t \in [0, 1]} \|f_{\theta^*(t)}(x) - f_{\gamma(t)}(x)\| \leq \frac{L_f \kappa}{8}. \quad (9)$$

¹¹

Nganjimi et al., Improving Clinical Dataset Condensation with Mode Connectivity-based Trajectory Surrogates, arXiv 2025.

Mode Connections as Trajectory Surrogates: Results



(a) Oxford University Hospitals (OUH)

Method	AUROC				AUPRC			
	50	100	200	500	50	100	200	500
RANDOM	0.835 _{+0.021}	0.855 _{+0.007}	0.869 _{+0.006}	0.888 _{+0.004}	0.158 _{+0.036}	0.176 _{+0.029}	0.219 _{+0.024}	0.276 _{+0.031}
M3D	0.840 _{+0.004}	0.862 _{+0.003}	0.872 _{+0.003}	0.893 _{+0.002}	0.162 _{+0.015}	0.190 _{+0.010}	0.266 _{+0.010}	0.290 _{+0.012}
MTT	0.824 _{+0.016}	0.849 _{+0.011}	0.855 _{+0.008}	0.870 _{+0.005}	0.356 _{+0.019}	0.381 _{+0.019}	0.405 _{+0.012}	0.407 _{+0.007}
TESLA	0.839 _{+0.007}	0.875 _{+0.002}	0.874 _{+0.003}	0.880 _{+0.002}	0.169 _{+0.012}	0.205 _{+0.009}	0.202 _{+0.010}	0.217 _{+0.011}
FTD	0.831 _{+0.014}	0.847 _{+0.006}	0.852 _{+0.005}	0.860 _{+0.006}	0.321 _{+0.031}	0.382 _{+0.012}	0.400 _{+0.009}	0.400 _{+0.023}
DATM	0.829 _{+0.010}	0.844 _{+0.007}	0.851 _{+0.009}	0.872 _{+0.009}	0.338 _{+0.020}	0.394 _{+0.009}	0.414 _{+0.009}	0.409 _{+0.004}
BTM (Ours)	0.854 _{+0.012}	0.863 _{+0.009}	0.876 _{+0.00}	0.888 _{+0.003}	0.394 _{+0.014}	0.396 _{+0.010}	0.427 _{+0.006}	0.436 _{+0.003}
<i>Full Dataset</i>	0.901_{+0.001}				0.445_{+0.004}			

(b) Portsmouth University Hospitals (PUH)

Method	AUROC				AUPRC			
	50	100	200	500	50	100	200	500
RANDOM	0.858 _{+0.006}	0.867 _{+0.004}	0.879 _{+0.004}	0.893 _{+0.004}	0.322 _{+0.025}	0.347 _{+0.011}	0.410 _{+0.029}	0.473 _{+0.018}
M3D	0.863 _{+0.004}	0.883 _{+0.003}	0.888 _{+0.002}	0.900 _{+0.001}	0.348 _{+0.013}	0.405 _{+0.011}	0.416 _{+0.011}	0.527 _{+0.011}
MTT	0.845 _{+0.009}	0.871 _{+0.005}	0.887 _{+0.005}	0.895 _{+0.004}	0.495 _{+0.013}	0.517 _{+0.012}	0.529 _{+0.013}	0.544 _{+0.018}
TESLA	0.861 _{+0.006}	0.861 _{+0.006}	0.885 _{+0.004}	0.893 _{+0.002}	0.412 _{+0.019}	0.385 _{+0.016}	0.472 _{+0.009}	0.483 _{+0.009}
FTD	0.845 _{+0.007}	0.870 _{+0.007}	0.885 _{+0.006}	0.899 _{+0.003}	0.499 _{+0.013}	0.536 _{+0.010}	0.550 _{+0.010}	0.584 _{+0.007}
DATM	0.871 _{+0.005}	0.876 _{+0.004}	0.887 _{+0.004}	0.889 _{+0.003}	0.540 _{+0.007}	0.555 _{+0.008}	0.555 _{+0.008}	0.579 _{+0.005}
BTM (Ours)	0.881 _{+0.007}	0.895 _{+0.007}	0.902 _{+0.003}	0.905 _{+0.002}	0.564 _{+0.011}	0.575 _{+0.008}	0.593 _{+0.009}	0.603 _{+0.007}
<i>Full Dataset</i>	0.906_{+0.002}				0.610_{+0.004}			

Mode Connections as Trajectory Surrogates: Results



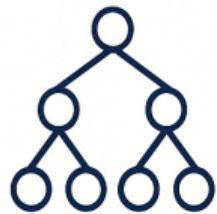
(c) University Hospitals Birmingham (UHB)

Method	AUROC				AUPRC			
	50	100	200	500	50	100	200	500
RANDOM	0.847 _{+0.015}	0.856_{+0.010}	0.876 _{+0.007}	0.891 _{+0.005}	0.089 _{+0.024}	0.108 _{+0.018}	0.126 _{+0.019}	0.153 _{+0.015}
M3D	0.863_{+0.005}	0.853 _{+0.003}	0.872 _{+0.004}	0.891_{+0.003}	0.107 _{+0.005}	0.070 _{+0.006}	0.107 _{+0.008}	0.153 _{+0.006}
MTT	0.802 _{+0.008}	0.847 _{+0.013}	0.871 _{+0.013}	0.884 _{+0.007}	0.092 _{+0.065}	0.228 _{+0.021}	0.242 _{+0.016}	0.234 _{+0.017}
TESLA	0.820 _{+0.001}	0.850 _{+0.004}	0.859 _{+0.004}	0.873 _{+0.006}	0.099 _{+0.001}	0.079 _{+0.006}	0.125 _{+0.008}	0.131 _{+0.014}
FTD	0.830 _{+0.005}	0.839 _{+0.010}	0.847 _{+0.010}	0.872 _{+0.010}	0.122 _{+0.056}	0.218 _{+0.013}	0.216 _{+0.018}	0.233 _{+0.012}
DATM	0.828 _{+0.028}	0.832 _{+0.022}	0.843 _{+0.015}	0.870 _{+0.005}	0.152 _{+0.019}	0.176 _{+0.022}	0.221 _{+0.017}	0.237 _{+0.017}
BTM (OURS)	0.842 _{+0.011}	0.836 _{+0.023}	0.884_{+0.007}	0.891 _{+0.006}	0.258_{+0.006}	0.246_{+0.013}	0.253_{+0.011}	0.272_{+0.006}
<i>Full Dataset</i>		0.895_{+0.005}			0.284_{+0.005}			

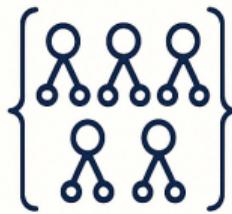
Dataset Condensation for Classical Clinical Models



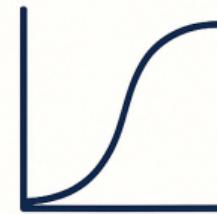
- ▶ **Classical AI models:** State-of-the-art (still!), regulatory familiarity and interpretability
- ▶ **Non-differentiable:** Existing DC methods require models to be differentiable
- ▶ Model-specific **inductive biases** can hinder condensed data **portability**



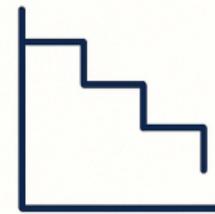
Decision
Tree



Random
Forest



Logistic
Regression



Cox
Model

Zero-order DC

- ▶ Real clinical data $\mathcal{X}_{\text{real}} = (\mathbf{X}, \mathbf{y})$, Synthetic dataset $\mathcal{X}_{\text{syn}} = (\mathbf{X}_{\text{syn}}, \mathbf{y}_{\text{syn}})$

- ▶ **Train a model on real data:**

$f_{\text{real}}(\cdot)$ (e.g., XGBoost) trained on $\mathcal{X}_{\text{real}}$

- ▶ **Initialise synthetic data:**

$\mathbf{X}_{\text{syn}} \in \mathbb{R}^{m \times d}$, $\mathbf{y}_{\text{syn}} \in \{0, 1\}^m$

- ▶ **Composite loss:**

$$\ell = \ell_{\text{pred}} + \alpha \cdot \ell_{\text{match}}$$

$$\ell_{\text{pred}} = \frac{1}{m} \sum_{i=1}^m [-y_i \log f_{\text{real}}(\mathbf{x}_i) - (1 - y_i) \log (1 - f_{\text{real}}(\mathbf{x}_i))]$$

$$\ell_{\text{match}} = \sum_{c \in \{0, 1\}} \left| \frac{1}{|\mathcal{I}_c^{\text{syn}}|} \sum_{i \in \mathcal{I}_c^{\text{syn}}} f_{\text{real}}(\mathbf{x}_i^{\text{syn}}) - \frac{1}{|\mathcal{I}_c^{\text{real}}|} \sum_{j \in \mathcal{I}_c^{\text{real}}} f_{\text{real}}(\mathbf{x}_j^{\text{real}}) \right|$$

- ▶ **Zero-order gradient:**

$$\nabla_{\mathbf{X}_{\text{syn}}} \ell = \left(\frac{\partial \ell}{\partial f_{\text{real}}(\mathbf{X}_{\text{syn}})} \right) \cdot \left(\frac{\partial f_{\text{real}}(\mathbf{X}_{\text{syn}})}{\partial \mathbf{X}_{\text{syn}}} \right)$$

- ▶ **Finite-diff approximation:**

$$\frac{\partial f_{\text{real}}(\mathbf{X}_{\text{syn}})}{\partial \mathbf{X}_{\text{syn},j}} \approx \frac{f_{\text{real}}(\mathbf{X}_{\text{syn}} + \epsilon_j \mathbf{E}_j) - f_{\text{real}}(\mathbf{X}_{\text{syn}} - \epsilon_j \mathbf{E}_j)}{2\epsilon_j}$$

- ▶ **Update rule:**

$$\mathbf{X}_{\text{syn}}^{(t+1)} = \mathbf{X}_{\text{syn}}^{(t)} - \eta \nabla_{\mathbf{X}_{\text{syn}}} \ell$$

Zero-order DC: Results



- **Datasets:** CURIAL datasets (PUH, OUH) and Proteomics
- **Baseline:** Real dataset

Portsmouth

INSTANCES	METRICS					PRIVACY BUDGET (ϵ)
	AUROC	SENSITIVITY	SPECIFICITY	PPV	NPV	
50	0.862 \pm 0.023	0.733 \pm 0.042	0.839 \pm 0.009	0.205 \pm 0.013	0.982 \pm 0.003	2.5
100	0.894 \pm 0.021	0.756 \pm 0.044	0.887 \pm 0.008	0.286 \pm 0.018	0.985 \pm 0.003	2.6
500	0.888 \pm 0.022	0.782 \pm 0.042	0.871 \pm 0.008	0.271 \pm 0.017	0.986 \pm 0.003	1.8
1000	0.884 \pm 0.022	0.741 \pm 0.041	0.868 \pm 0.007	0.262 \pm 0.016	0.983 \pm 0.003	3.1
<i>Full Dataset</i>	0.901 \pm 0.022	0.751 \pm 0.045	0.905 \pm 0.006	0.316 \pm 0.026	0.985 \pm 0.003	–

Oxford

INSTANCES	METRICS					PRIVACY BUDGET (ϵ)
	AUROC	SENSITIVITY	SPECIFICITY	PPV	NPV	
50	0.855 \pm 0.019	0.752 \pm 0.037	0.819 \pm 0.007	0.067 \pm 0.004	0.995 \pm 0.001	2.2
100	0.873 \pm 0.019	0.744 \pm 0.038	0.847 \pm 0.003	0.077 \pm 0.004	0.995 \pm 0.001	2.7
500	0.870 \pm 0.019	0.725 \pm 0.039	0.843 \pm 0.004	0.073 \pm 0.004	0.994 \pm 0.001	2.1
1000	0.891 \pm 0.018	0.760 \pm 0.038	0.887 \pm 0.005	0.095 \pm 0.005	0.995 \pm 0.001	1.9
<i>Full Dataset</i>	0.911 \pm 0.018	0.756 \pm 0.035	0.898 \pm 0.003	0.126 \pm 0.008	0.996 \pm 0.001	–

Proteomics

INSTANCES	METRICS					PRIVACY BUDGET (ϵ)
	AUROC	SENSITIVITY	SPECIFICITY	PPV	NPV	
50	0.856 \pm 0.053	0.788 \pm 0.106	0.817 \pm 0.003	0.014 \pm 0.002	0.999 \pm 0.001	2.1
100	0.905 \pm 0.048	0.846 \pm 0.093	0.785 \pm 0.006	0.013 \pm 0.001	0.999 \pm 0.001	2.6
500	0.913 \pm 0.044	0.865 \pm 0.092	0.834 \pm 0.006	0.017 \pm 0.002	0.999 \pm 0.001	1.9
1000	0.912 \pm 0.052	0.808 \pm 0.096	0.879 \pm 0.005	0.022 \pm 0.003	0.999 \pm 0.001	2.3
<i>Full Dataset</i>	0.898 \pm 0.050	0.846 \pm 0.091	0.757 \pm 0.006	0.011 \pm 0.001	0.999 \pm 0.001	–

Zero-order DC: External Validation using CURIAL



MODEL	PUH → UHB					OUH → UHB				
	AUROC	SENSITIVITY	SPECIFICITY	PPV	NPV	AUROC	SENSITIVITY	SPECIFICITY	PPV	NPV
FULL	0.830 ± 0.037	0.631 ± 0.074	0.866 ± 0.005	0.038 ± 0.005	0.996 ± 0.001	0.842 ± 0.040	0.611 ± 0.077	0.917 ± 0.004	0.058 ± 0.008	0.996 ± 0.001
IPC 50	0.855 ± 0.033	0.758 ± 0.064	0.847 ± 0.006	0.040 ± 0.004	0.998 ± 0.001	0.825 ± 0.037	0.711 ± 0.074	0.814 ± 0.006	0.031 ± 0.003	0.997 ± 0.001
IPC 100	0.860 ± 0.035	0.765 ± 0.070	0.840 ± 0.004	0.039 ± 0.003	0.998 ± 0.001	0.831 ± 0.035	0.691 ± 0.074	0.863 ± 0.005	0.041 ± 0.004	0.997 ± 0.001
IPC 500	0.847 ± 0.038	0.698 ± 0.074	0.836 ± 0.003	0.035 ± 0.004	0.997 ± 0.001	0.864 ± 0.033	0.772 ± 0.064	0.833 ± 0.006	0.037 ± 0.003	0.998 ± 0.001
IPC 1000	0.861 ± 0.032	0.805 ± 0.064	0.845 ± 0.005	0.042 ± 0.003	0.998 ± 0.001	0.844 ± 0.035	0.745 ± 0.070	0.862 ± 0.005	0.043 ± 0.004	0.998 ± 0.001

MODEL	PUH → OUH					UHB → OUH				
	AUROC	SENSITIVITY	SPECIFICITY	PPV	NPV	AUROC	SENSITIVITY	SPECIFICITY	PPV	NPV
FULL	0.863 ± 0.021	0.705 ± 0.041	0.885 ± 0.003	0.095 ± 0.006	0.994 ± 0.001	0.778 ± 0.020	0.326 ± 0.041	0.953 ± 0.002	0.085 ± 0.013	0.988 ± 0.001
IPC 60	0.862 ± 0.019	0.723 ± 0.039	0.844 ± 0.004	0.074 ± 0.004	0.994 ± 0.001	0.806 ± 0.022	0.502 ± 0.043	0.891 ± 0.003	0.070 ± 0.006	0.991 ± 0.001
IPC 100	0.873 ± 0.019	0.721 ± 0.039	0.874 ± 0.004	0.089 ± 0.005	0.995 ± 0.001	0.821 ± 0.021	0.613 ± 0.040	0.870 ± 0.004	0.072 ± 0.005	0.993 ± 0.001
IPC 500	0.848 ± 0.020	0.727 ± 0.039	0.823 ± 0.004	0.066 ± 0.004	0.994 ± 0.001	0.806 ± 0.022	0.525 ± 0.041	0.886 ± 0.003	0.070 ± 0.005	0.991 ± 0.001
IPC 1000	0.851 ± 0.020	0.729 ± 0.039	0.848 ± 0.004	0.076 ± 0.004	0.995 ± 0.001	0.803 ± 0.022	0.517 ± 0.042	0.871 ± 0.004	0.062 ± 0.005	0.991 ± 0.001

Table of Contents

Preliminary Concepts

Transforms

Loss Landscape & Mode Connections

Multi-tasking

Adaptive Parameter Optimisation

Task Groupings

Gradient-based Meta-Learning

Federated Learning

Data Democratisation & Healthcare

Data Encoding and Condensation

Uncertainty in LLMs



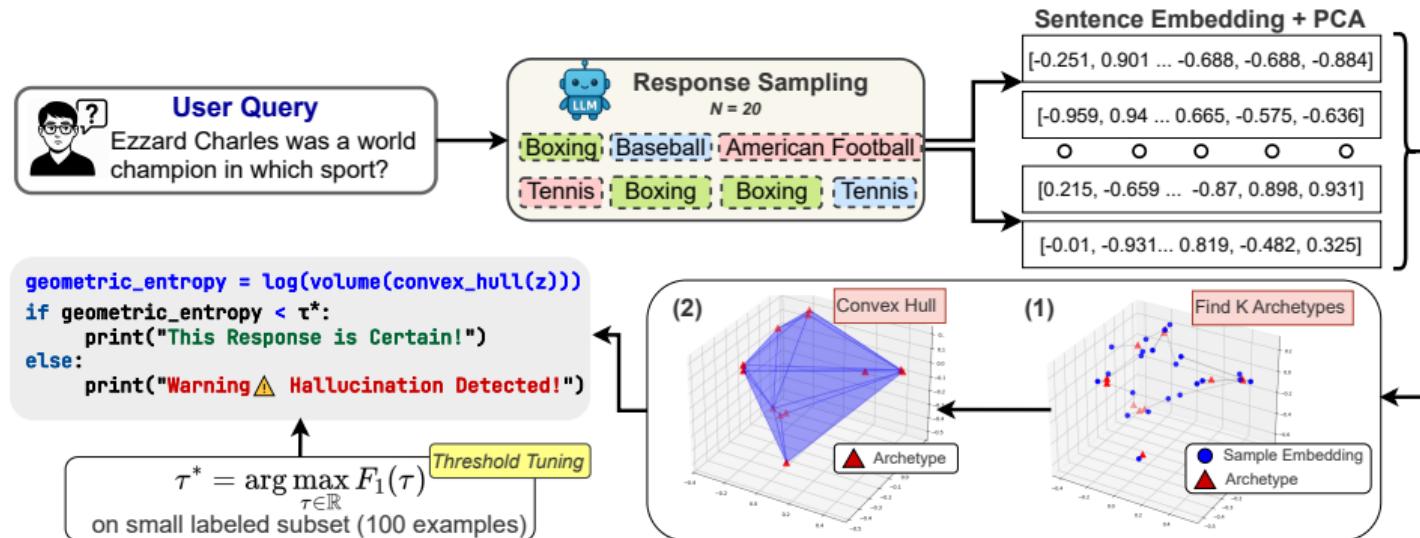
Uncertainty in LLMs and Its Link to Hallucinations



- ▶ LLMs are probabilistic systems that generate outputs based on learned distributions
- ▶ **Uncertainty** arises when the model has low confidence in its output or faces ambiguous input
- ▶ Understanding uncertainty is crucial to detecting hallucinations – responses that are incorrect, fabricated, or misleading.

Type	Description	Example	Link to Hallucination
External Uncertainty	Ambiguity or underspecified user intent	Who played Spiderman?	Model guesses instead of clarifying intent
Internal Uncertainty	Missing, outdated, or conflicting model knowledge	What's the latest COVID variant?	Model fills gap with hallucinated facts

Simplicial Modelling-Based Uncertainty Estimation¹²



¹²

Philips et al., Geometric Uncertainty for Detecting and Correcting Hallucinations in LLMs, arXiv 2025.

Benchmarks and Results: Global Uncertainty

Experimental Setup

Benchmark	LLM(s)
CLAMBER	GPT-4o Mini
TriviaQA	LLaMA-3.2B Instruct
K-QA	GPT-3.5 Turbo, GPT-4o Mini
Med-QA	GPT-4o Mini
Sentence Encoder:	gte-Qwen2-1.5B-Instruct

AUROC and F1 Score Comparison

AUROC

Dataset	p(true)	Semantic Entropy	Semantic Volume	Ours
CLAMBER	58.5 _{0.4}	48.1 _{0.8}	55.8 _{0.9}	61.7 _{0.7}
TriviaQA	73.3 _{5.6}	79.3 _{0.2}	75.5 _{0.3}	75.6 _{0.4}
K-QA(3.5)	60.5 _{1.2}	66.4 _{1.2}	65.3 _{0.7}	67.6 _{1.2}
K-QA(4o)	63.4 _{1.5}	50.5 _{5.7}	65.2 _{1.1}	69.7 _{2.3}
Med-QA(4o)	59.4 _{0.9}	59.3 _{1.1}	59.4 _{1.2}	62.3 _{1.3}

F1 Score

Dataset	p(true)	Semantic Entropy	Semantic Volume	Ours
CLAMBER	47.0 _{1.6}	63.9 _{0.3}	68.5 _{0.4}	66.1 _{0.3}
TriviaQA	69.4 _{3.0}	74.9 _{0.3}	73.4 _{1.6}	74.7 _{0.3}
K-QA(3.5)	71.8 _{1.4}	40.3 _{6.9}	74.8 _{1.0}	75.6 _{0.9}
K-QA(4o)	61.5 _{4.7}	70.4 _{4.8}	78.2 _{1.6}	79.2 _{1.4}
Med-QA(4o)	73.6 _{1.2}	70.9 _{1.1}	72.6 _{1.2}	73.1 _{1.4}

Note: **Bold** = highest per row. Subscripts = standard deviation.

Thank You!

