



pFedNavi: Structure-Aware Personalized Federated Vision-Language Navigation for Embodied AI

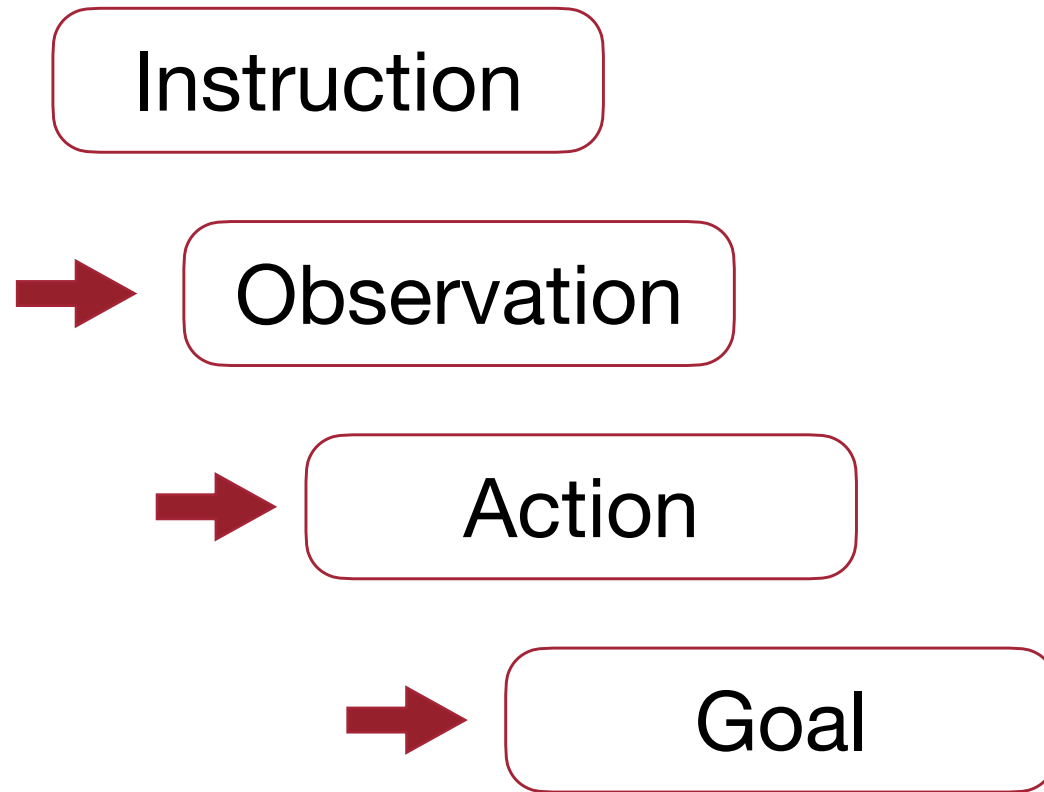
Qingqian Yang¹, Hao Wang¹, Sai Qian Zhang²,
Jian Li³, Yang Hua⁴, Miao Pan⁵, Tao Song⁶,
Zhengwei Qi⁶, and Haibin Guan⁶

 IntelliSys Lab

Stevens Institute of Technology¹, New York University²,
Stony Brook University³, Queen's University Belfast⁴,
University of Houston⁵, Shanghai Jiao Tong University⁶

Vision Language Navigation (VLN)

An embodied agent follows natural-language instructions to reach a target location.



Client i : “Fetch a bottle of drink from kitchen to bedroom”



Example of VLN task

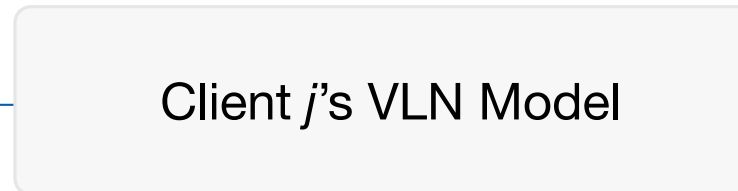
Federated VLN: Privacy-Preserving Training



Client i : “Fetch a bottle of drink from kitchen to bedroom”



-
-
-

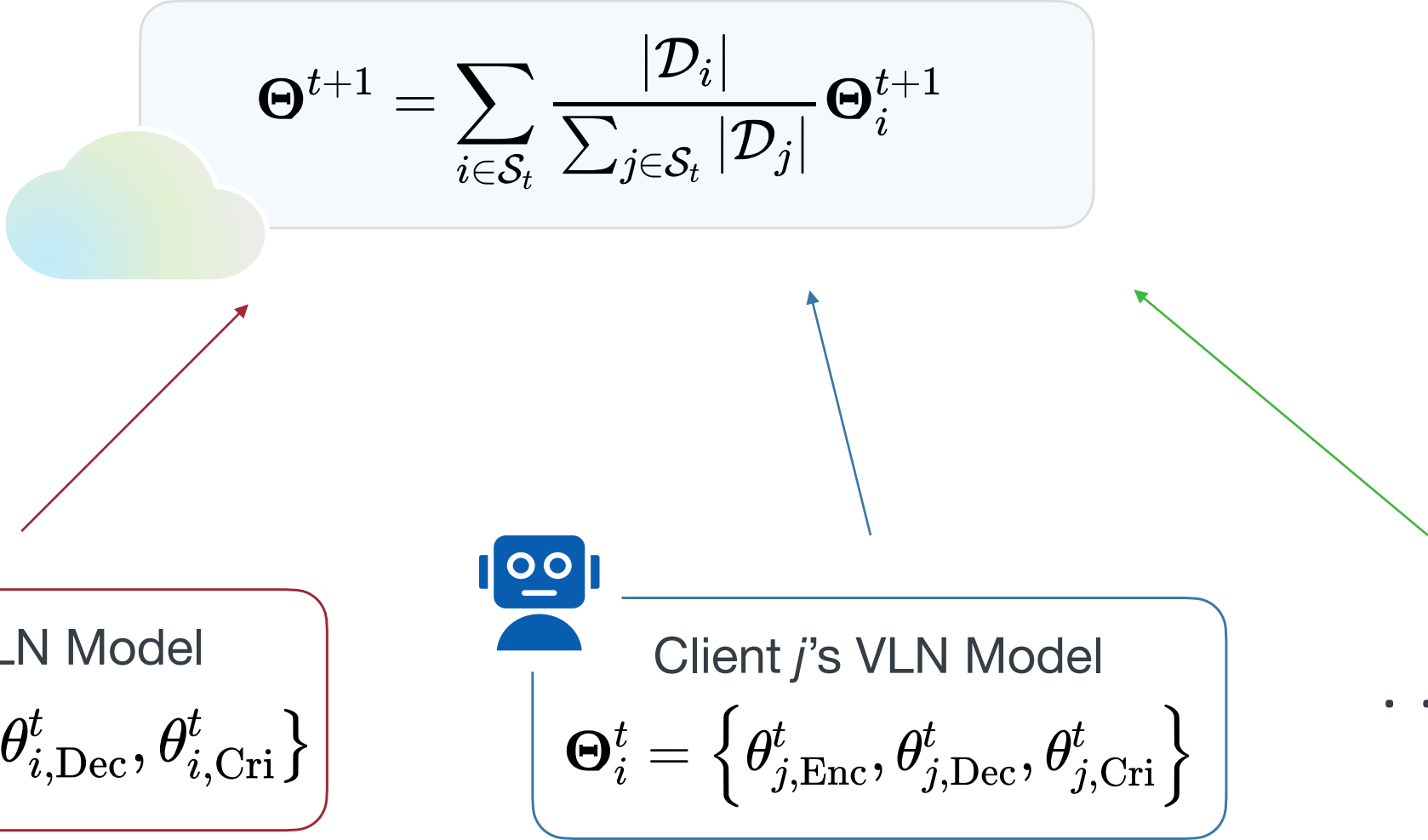


Local upload ↗ ↘ *Aggregate*

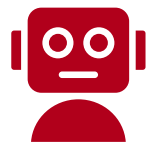


Federated Server

Existing Works on Federated VLN (*FedVLN*)

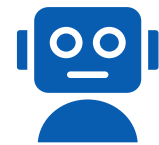


The diagram illustrates the Federated VLN (FedVLN) architecture. At the top, a light blue rounded rectangle contains the aggregation formula:
$$\Theta^{t+1} = \sum_{i \in \mathcal{S}_t} \frac{|\mathcal{D}_i|}{\sum_{j \in \mathcal{S}_t} |\mathcal{D}_j|} \Theta_i^{t+1}$$
 To the left of this box is a stylized cloud icon. Below the cloud, a red arrow points from a red robot icon (representing Client i) to the cloud. A blue arrow points from a blue robot icon (representing Client j) to the cloud. A green arrow points from the cloud to the right. Below the cloud, there are two rounded rectangles representing client models. The left one is red and labeled 'Client i 's VLN Model', containing the equation
$$\Theta_i^t = \{ \theta_{i,Enc}^t, \theta_{i,Dec}^t, \theta_{i,Cri}^t \}$$
 The right one is blue and labeled 'Client j 's VLN Model', containing the equation
$$\Theta_j^t = \{ \theta_{j,Enc}^t, \theta_{j,Dec}^t, \theta_{j,Cri}^t \}$$
 To the right of the blue model box are three dots '...'.



Client i 's VLN Model

$$\Theta_i^t = \{ \theta_{i,Enc}^t, \theta_{i,Dec}^t, \theta_{i,Cri}^t \}$$



Client j 's VLN Model

$$\Theta_j^t = \{ \theta_{j,Enc}^t, \theta_{j,Dec}^t, \theta_{j,Cri}^t \}$$

...

Key Problems — Heterogeneity

- **Personalized Instructions:**

- ▶ *different vocabulary*
- ▶ *detail levels*
- ▶ *landmarks*



Client *i*: “Fetch a bottle of drink from kitchen to bedroom”

Client *j*: “Go through the hallway ... to the sofa”

Client *k*: “Go to the balcony via the living room”



...



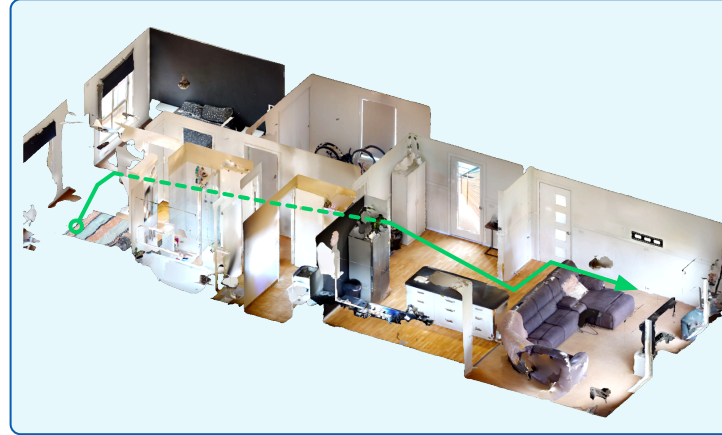
Key Problems — Heterogeneity



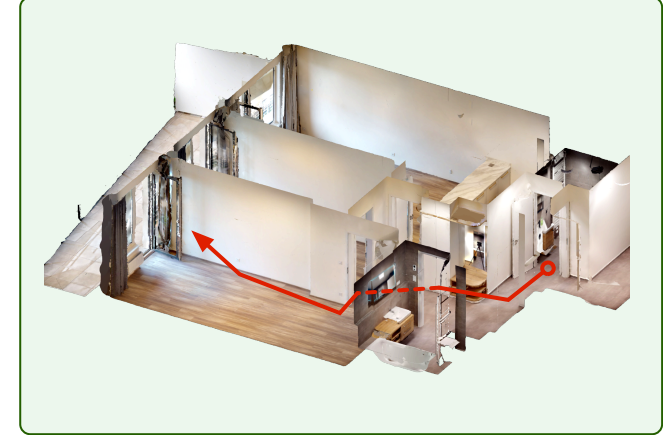
Client *i*: “Fetch a bottle of drink from kitchen to bedroom”



Client *j*: “Go through the hallway ... to the sofa”



Client *k*: “Go to the balcony via the living room”



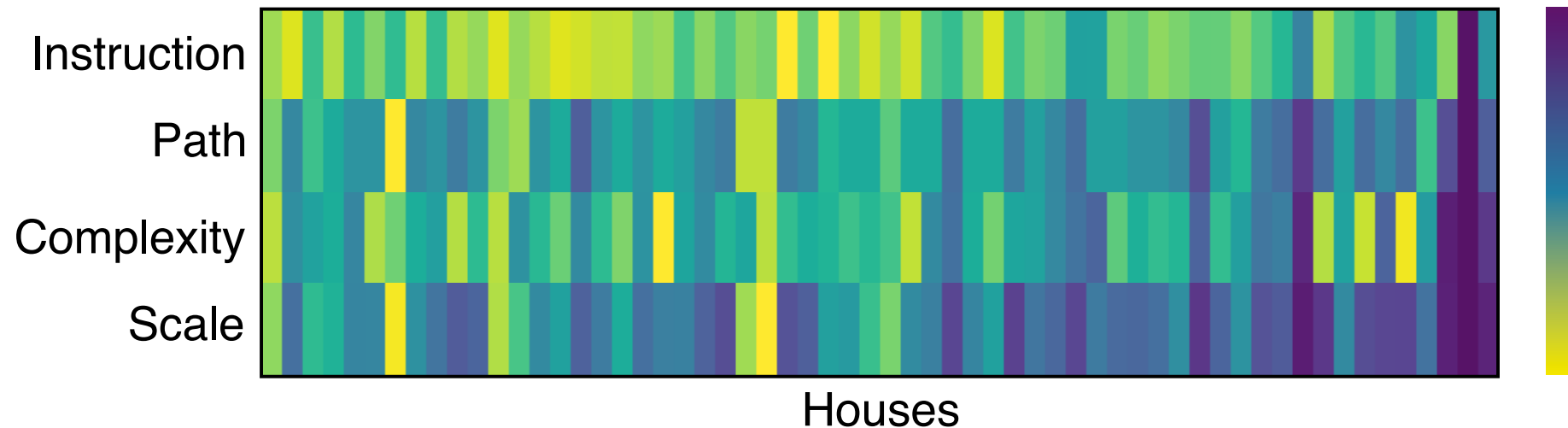
■ Heterogeneous Environments:

- ▶ *layout*
- ▶ *scale*

- ▶ *graph complexity*
- ▶ *path distribution*

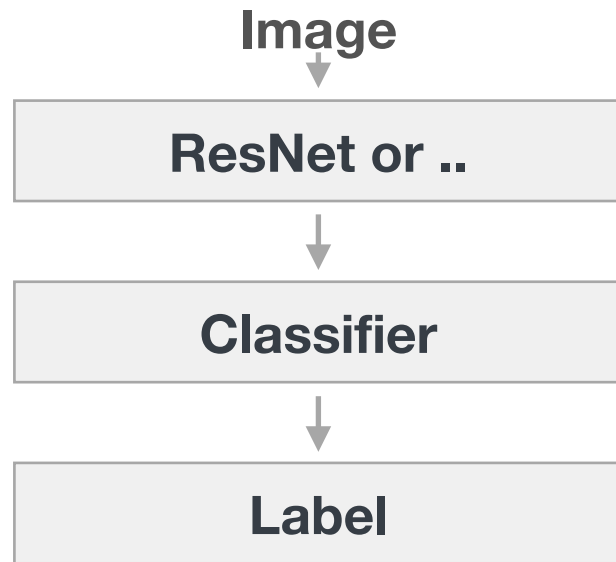
Key Problems — Heterogeneity

- Data heterogeneity analysis on RxR dataset [1]:
 - ▶ *layout*
 - ▶ *scale*
 - ▶ *graph complexity*
 - ▶ *path distribution*



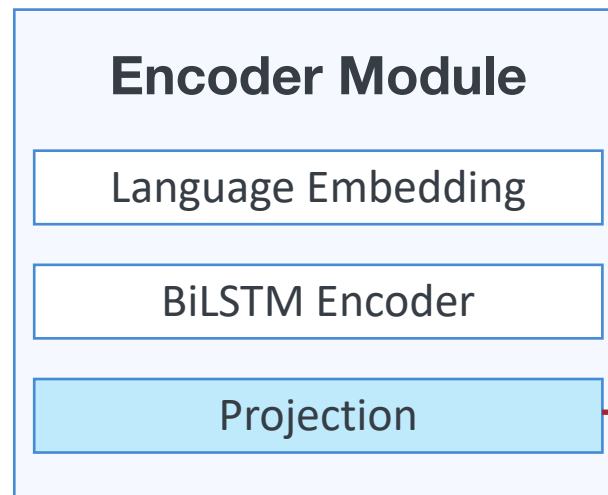
[1] Ku, Alexander, et al. "Room-across-room: Multilingual vision-and-language navigation with dense spatiotemporal grounding." *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*. 2020.

Key Problems – Multimodal Models

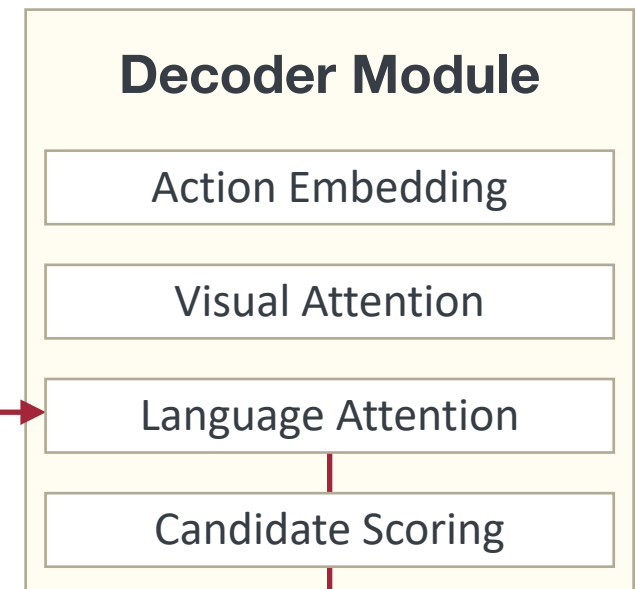


Traditional Vision Model

Language instructions



Visual observation



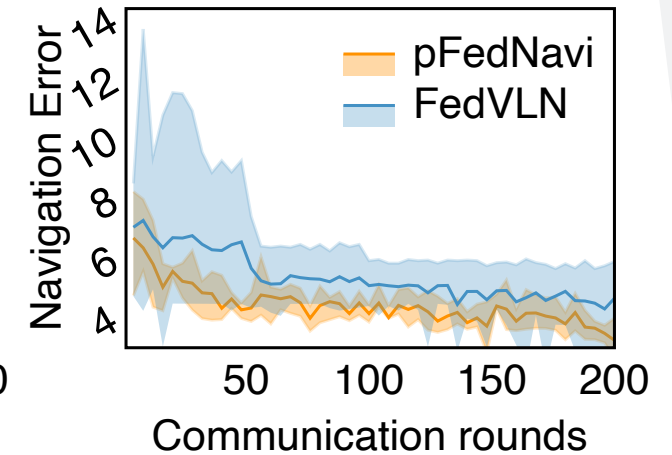
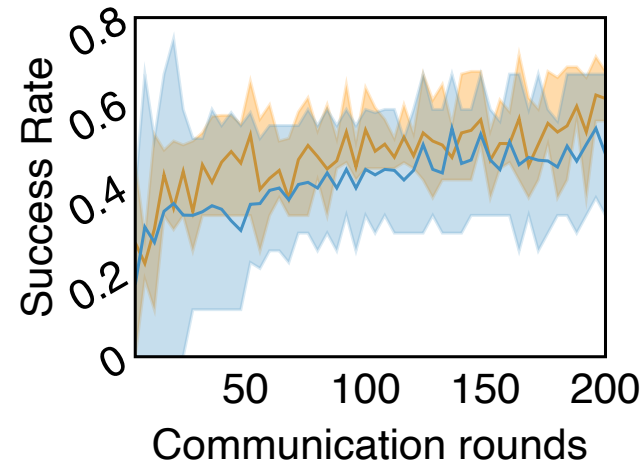
VLN Navigation model

Personalized Federated VLN

Heterogeneous
Environments

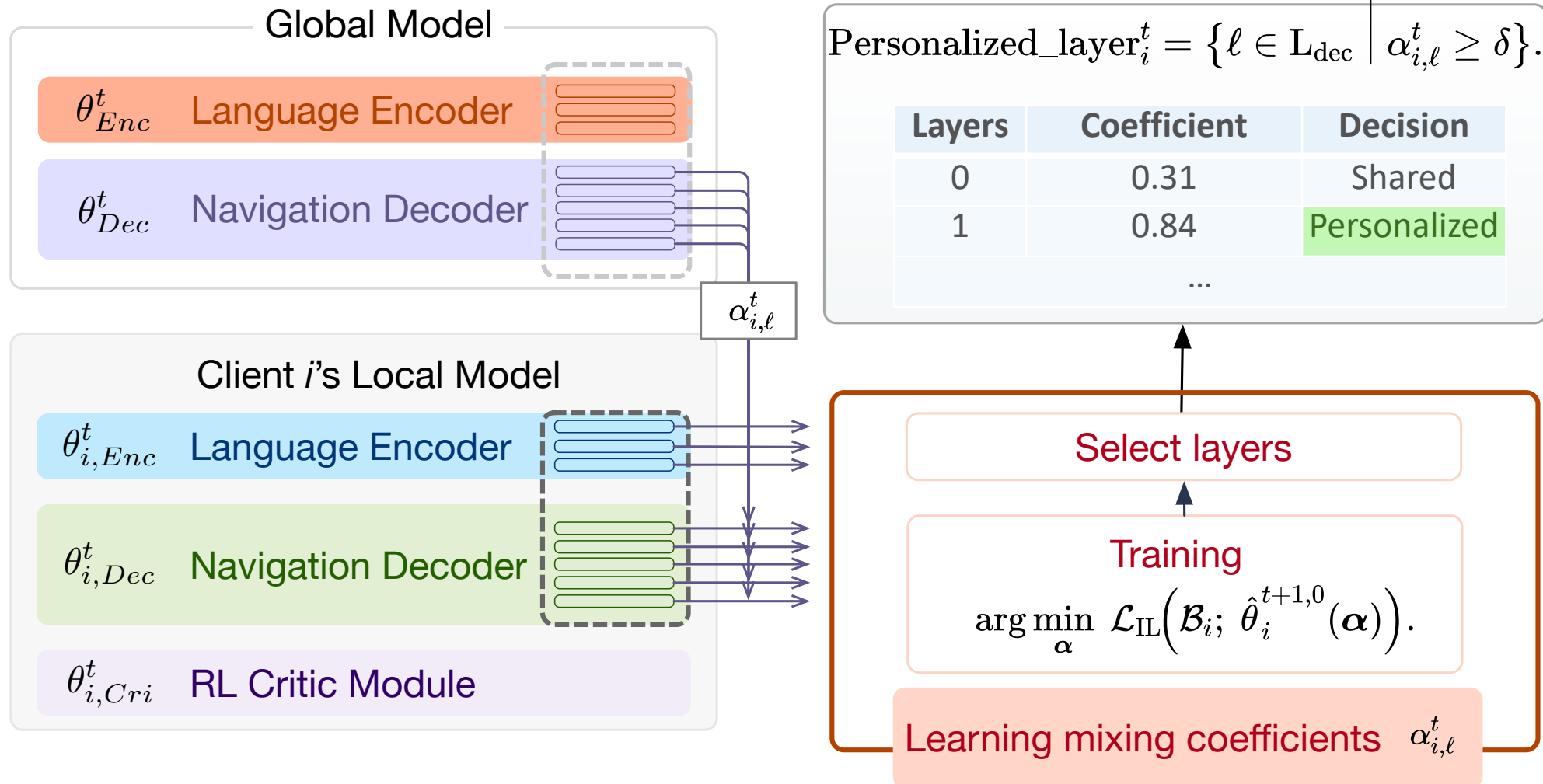
Personalized Instructions

Complex Multimodal
Models



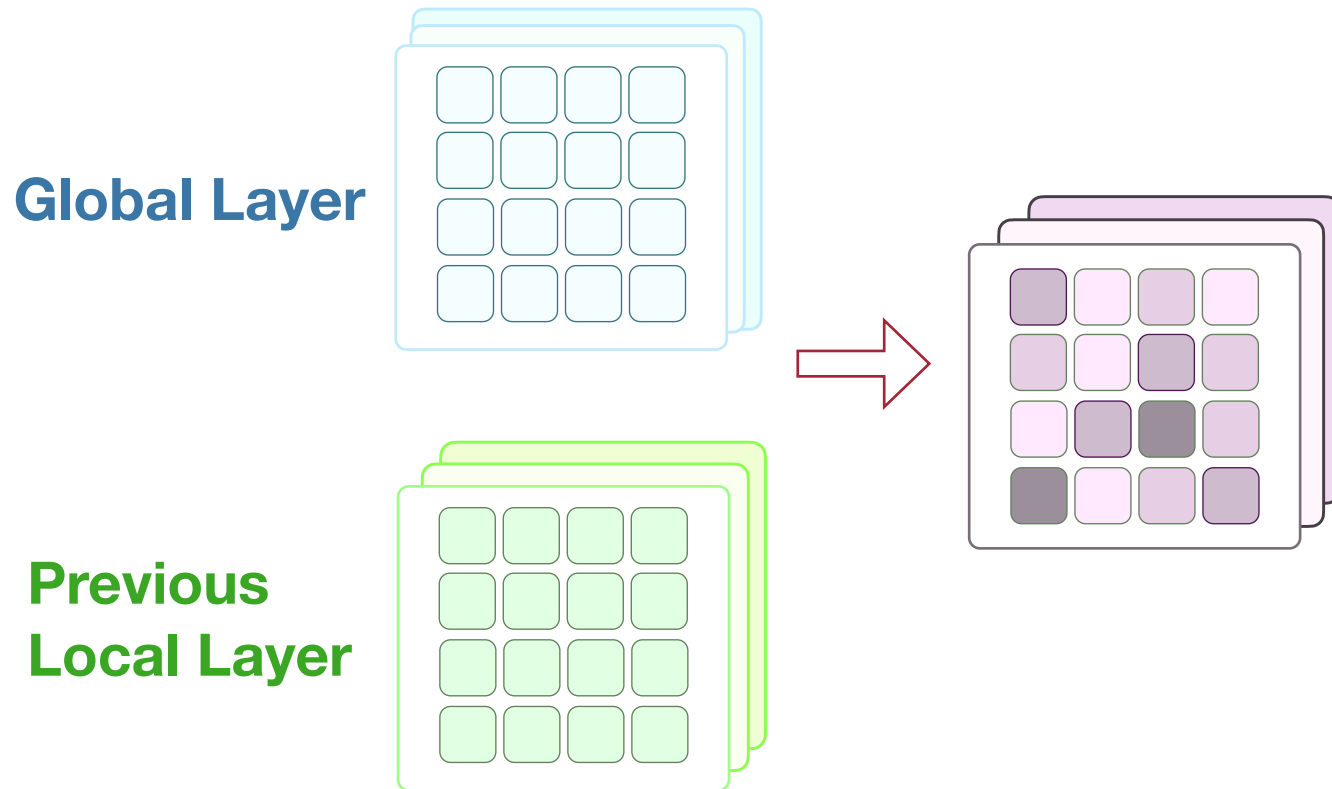
Non-personalized FL (FedVNL) vs. personalized FL (pFedNavi).

Adaptive Layer Selection for Personalization



Fine-grained Parameter Fusion

For each selected personalized layers:



Fusion formulation

$$\theta_{i,l}^{t+1,0} = \theta_{i,l}^t + \mathbf{W}_{i,l}^t \odot (\theta_l^t - \theta_{i,l}^t)$$

Optimized by:

$$\min_{\mathbf{W}_i^t} \mathcal{L}_{\text{IL}}(\mathcal{D}_i; \Theta_i^{t+1,0}(\mathbf{W}_i^t))$$

Evaluation

■ Datasets:

- ▶ R2R
- ▶ RxR

■ Models:

- ▶ Seq2seq model

■ Baselines:

- ▶ EnvDrop (centralized setting)
- ▶ FedVLN (federated setting)

■ Pretrained Image features:

- ▶ ResNet-152 representations
- ▶ CLIP representations

■ Metrics:

- ▶ Success Rate (SR)
- ▶ Success weighted by Path Length (SPL)
- ▶ Oracle Success Rate (OSR)
- ▶ Navigation Error (NE)
- ▶ Coverage weighted by Length Score (CLS)
- ▶ Normalized Dynamic Time Warping (nDTW)

Main Results

Improves Goal-Reaching Ability

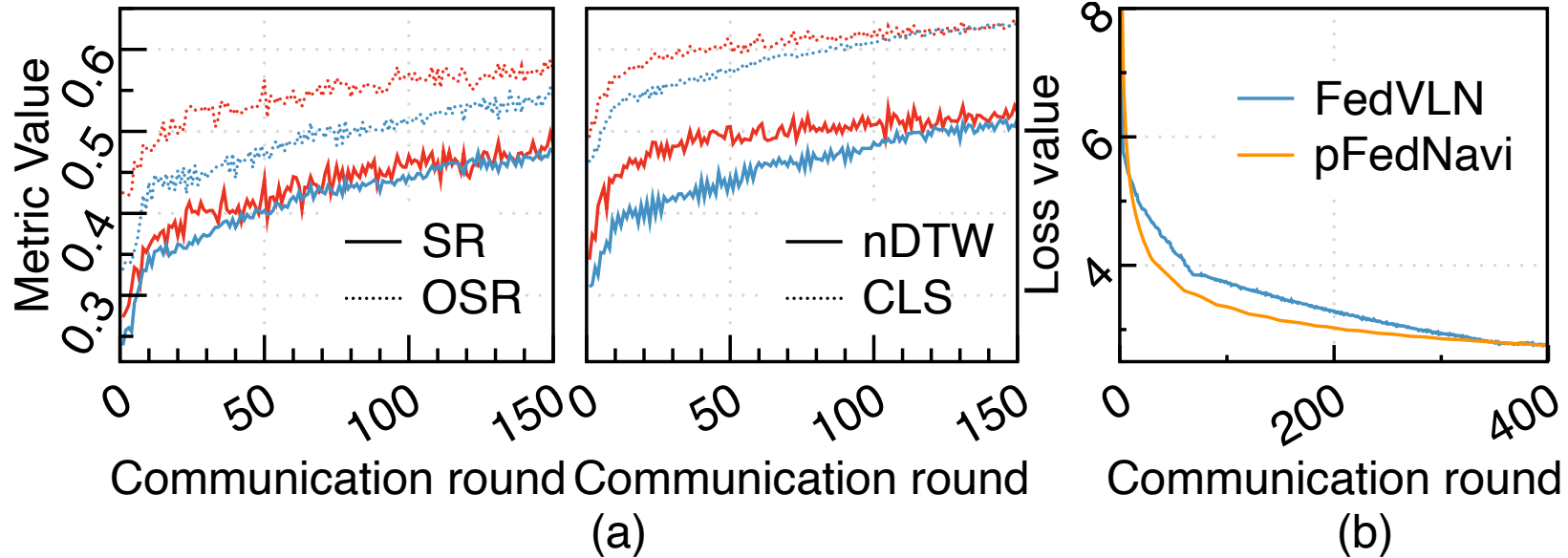
up to 7.5%

Improves Trajectory Fidelity

up to 7.8%

Method on R2R		SR ↑	SPL ↑	OSR ↑	CLS ↑	nDTW ↑	NE ↓
ResNet	EnvDrop [9]	56.7	54.3	63.9	67.2	56.0	4.55
	FedVLN [2]	50.7	47.3	60.3	63.4	50.8	5.37
	<i>pFedNavi</i>	54.5	51.7	65.2	66.4	54.8	4.86
Method on RxR		SR ↑	SPL ↑	OSR ↑	CLS ↑	nDTW ↑	NE ↓
ResNet	EnvDrop [9]	41.3	36.6	50.4	57.2	53.6	8.03
	FedVLN [2]	39.9	33.1	47.8	55.5	50.7	8.62
	<i>pFedNavi</i>	40.1	36.7	51.2	57.3	52.9	8.10
CLIP	EnvDrop [9]	47.7	43.4	56.1	61.2	56.8	6.53
	FedVLN [2]	43.0	39.4	51.5	58.7	54.8	7.37
	<i>pFedNavi</i>	46.1	41.2	56.3	59.8	56.9	6.91

Main Results



Higher efficiency: 1.38x faster convergence

Various personalized layer selection strategies

pFedNavi

All Layers

-18.25% on SPL

No Layers

-9.80% on SPL



IntelliSys Lab



STEVENS
INSTITUTE OF TECHNOLOGY
1870



Code



Paper

Stevens Institute of Technology
1 Castle Point Terrace, Hoboken, NJ 07030