

MLG







Matthieu Kirchmeyer\*<sup>1,2</sup>, Yuan Yin\*<sup>1</sup>, Jérémie Donà<sup>1</sup>, Nicolas Baskiotis<sup>1</sup>, Alain Rakotomamonjy<sup>2,3</sup>, Patrick Gallinari<sup>1,2</sup>

\*Equal Contribution

<sup>1</sup>Sorbonne Université, ISIR, MLIA <sup>2</sup>Criteo AI Lab

<sup>3</sup>Université de Rouen, LITIS

{matthieu.kirchmeyer,yuan.yin}@sorbonne-universite.fr





### MOTIVATION

- Neural dynamics models successful for modeling a physical system but fail on out-of-distribution systems.
- Limitation for real-world problems:
- predicting disease diffusion in new countries
- modelling heart blood flow for new patients
- predicting ocean dynamics for new spatial regions on earth
- CoDA (Context-Informed Dynamics Adaptation): first principled solution to this problem.

### PROBLEM SETTING

### DYNAMICS-AWARE FORMULATION

We consider dynamical systems described by a differential equation (ODE/PDE):

$$\frac{\mathrm{d}x(t)}{\mathrm{d}t} = f(x(t))$$

- x(t): state at t
- f: unknown dynamics describing the evolution of x
- depends on context e.g. parameters, forcing
- defines trajectories:  $x(t) = x_0 + \int_0^t f(x(\tau)) d\tau$

### LEARNING ACROSS ENVIRONMENTS

We learn a neural dynamics model  $g_{\theta}$  across contexts.

- We leverage several different environments:
- environment  $e \in \mathcal{E} \Leftrightarrow$  physical context
- trajectories  $\mathcal{D}^e$  of corresponding dynamics  $f^e$
- Training: environments  $\mathcal{E}_{\mathrm{tr}}$  with reasonable data
- Adaptation: environments  $\mathcal{E}_{\mathrm{ad}}$  with scarce data
- Task: accurately predict new trajectories of  $\mathcal{E}_{\mathrm{ad}}$

# THEORETICAL MOTIVATION

**Proposition 1** (Low-rank gradients). For linearly parameterized dynamics with  $d_p$  parameters,  $\forall \theta^c \in \mathbb{R}^{d_\theta}$ ,  $\dim(\operatorname{Span}(\{\nabla_{\theta}\mathcal{L}(\theta^c, \mathcal{D}^e)\}_{e \in \mathcal{E}})) \leq d_p \ll d_{\theta}$ .

ightharpoonup Gradients of MSE loss  $\mathcal{L}(\theta, \mathcal{D}^e)$  across environments live in a tiny subspace.

### CODA FRAMEWORK

### ADAPTATION RULE

$$\forall e, \theta^e \triangleq \theta^c + \delta \theta^e$$

 $heta^c$  shared;  $\delta heta^e$  environment-specific parameters  $\in \mathbb{R}^{d_ heta}$ 

#### LOCALITY

$$\min_{\theta^c, \{\delta\theta^e\}} \sum_{e \in \mathcal{E}} \lVert \delta\theta^e \rVert^2 \text{ s.t. } \forall t \ \frac{\mathrm{d} x^e(t)}{\mathrm{d} t} = g_{\theta^c + \delta\theta^e}(x^e(t))$$

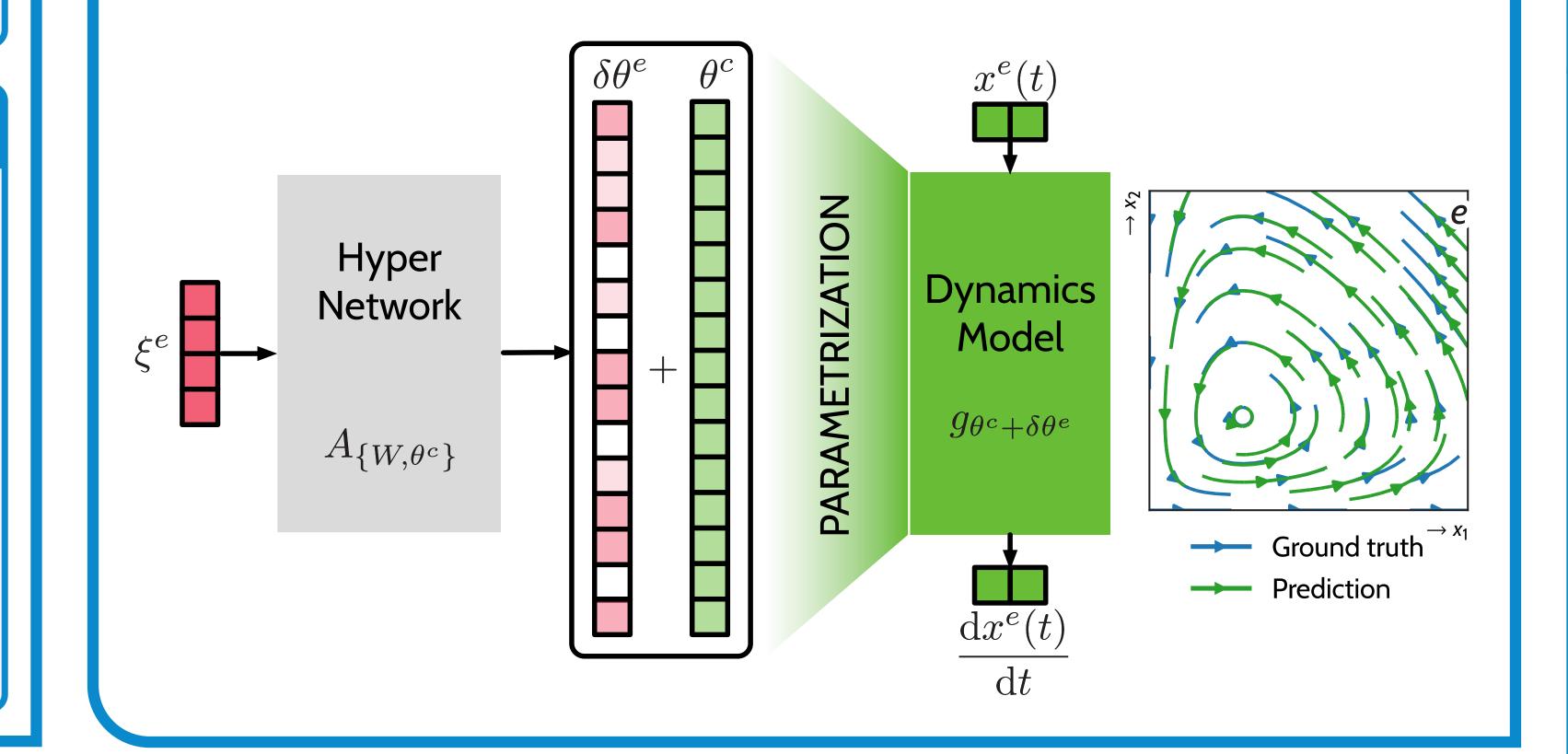
- Fast adaptation by constraining  $\theta^c$  during training  $\rightarrow$  few update steps
- Hypothesis space constrained around  $\theta^c$
- → under assumptions, optimization is quadratic and convex

#### LOW-RANK ADAPTATION

 $\theta^e$  generated via a trained hypernetwork:

$$\forall e, \theta^e \triangleq \theta^c + W\xi^e \quad (\delta\theta^e \triangleq W\xi^e)$$

- $\xi^e$ : context vector  $\in \mathbb{R}^{d_\xi}$  ( $d_\xi \ll d_ heta$ )
- → smaller adaptation space



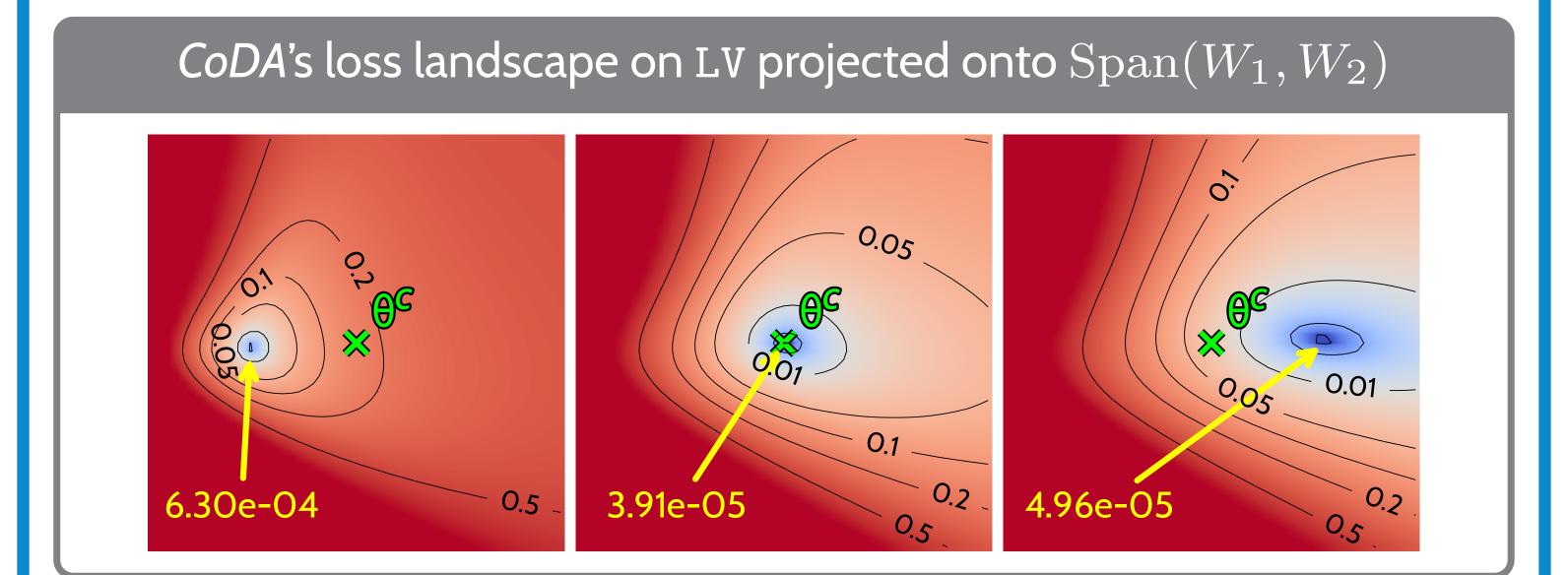
### BENEFITS

- Fast and sample-efficient adaptation
- Time-continuous
- Agnostic to the architecture of  $g_{\theta}$ :

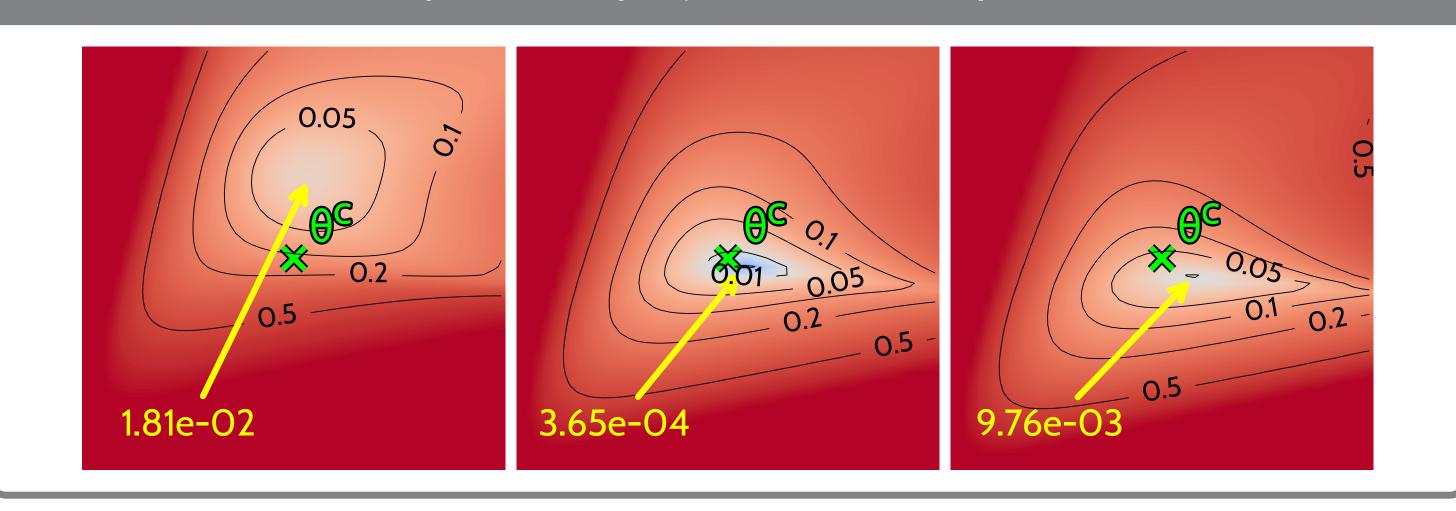
 $ODE \rightarrow MLP \quad PDE \rightarrow ConvNet, FNO$ 

 Efficient task conditioning: hypernet decoding generalizes FiLM / Concatenation conditioning approaches

### LOSS LANDSCAPES







- Smooth loss with single minimum across environments
- Proximity of local loss optimas to  $\theta^c$
- Lower minimal loss value for CoDA than ERM

### MPLEMENTATION

### OBJECTIVE FUNCTION

Training: 
$$\min_{\theta^c, W, \{\xi^e\}_{e \in \mathcal{E}_{\mathrm{tr}}}} \sum_{e \in \mathcal{E}_{\mathrm{tr}}} \mathcal{L}(\theta^c + W\xi^e, \mathcal{D}^e) + \|W\xi^e\|^2$$

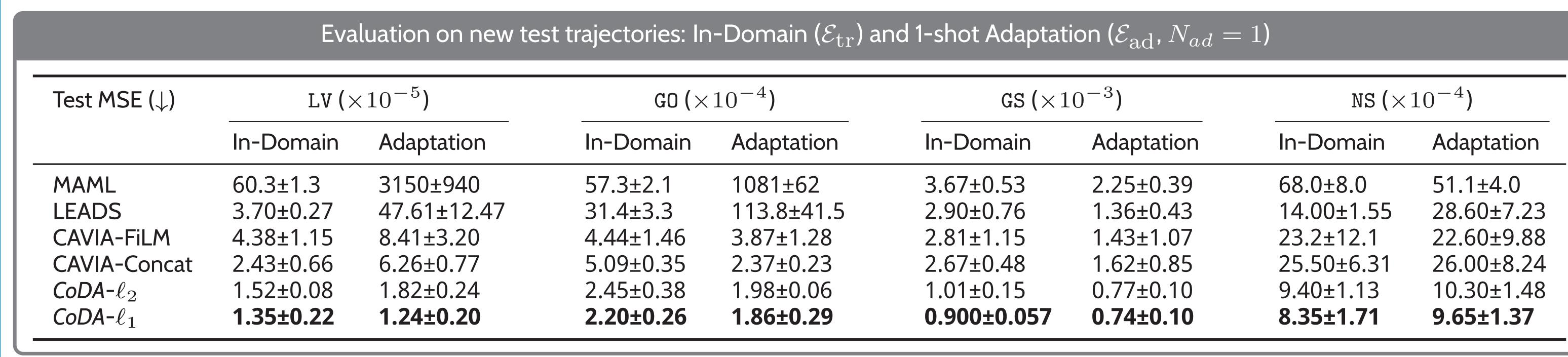
Adaptation:  $\min_{\{\xi^e\}_{e\in\mathcal{E}_{ad}}} \sum_{e\in\mathcal{E}_{ad}} \mathcal{L}(\theta^c + W\xi^e, \mathcal{D}^e) + \|W\xi^e\|^2$ 

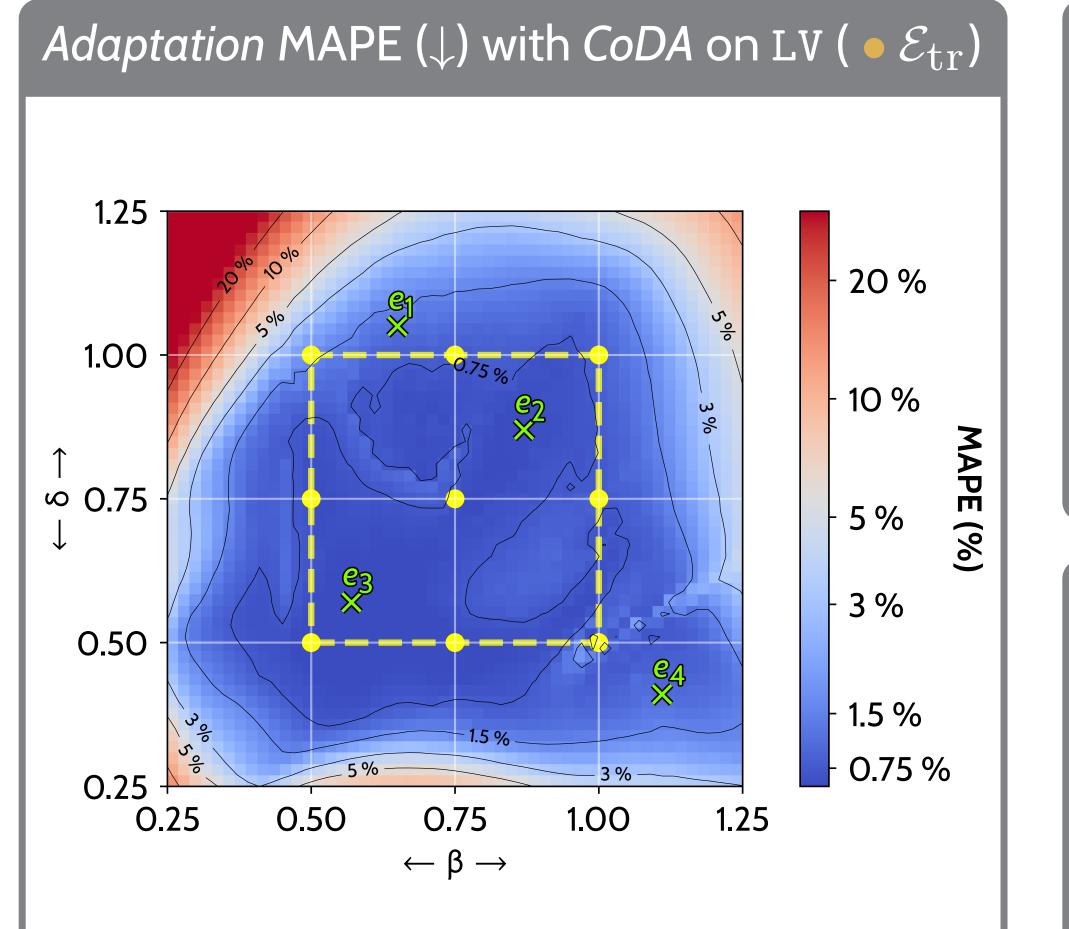
- Loss:  $\mathcal{L}(\theta, \mathcal{D}^e) = \sum_{i=1}^{N_{tr}} \sum_{t_k} \left\| (x^{e,i} \tilde{x}^{e,i})(t_k) \right\|_2^2$  where  $\tilde{x}^{e,i}(t_k) = x_0^{e,i} + \int_0^{t_k} g_{\theta}(\tilde{x}^{e,i}(\tau)) d\tau$
- Regularization:  $\|W\xi^e\|^2 \to \lambda_{\xi} \|\xi^e\|_2^2 + \lambda_{\Omega}\Omega(W)$
- $\ell_2$ :  $\Omega(W) = ||W||_2$   $\ell_1$ :  $\Omega(W) = \sum_{i=1}^{d_\theta} ||W_{i,\cdot}||_2$

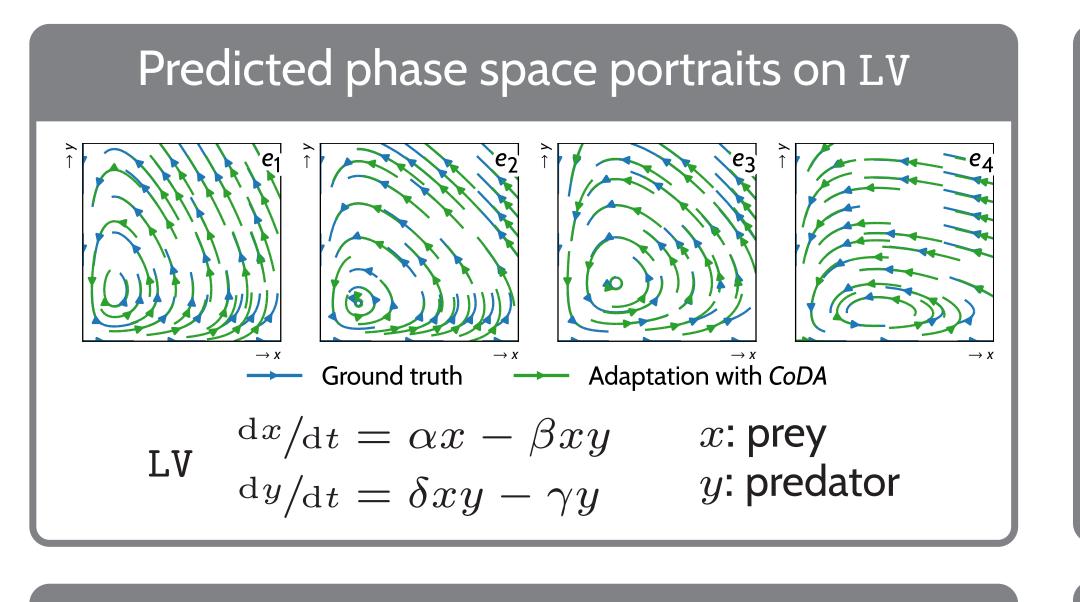
# EXPERIMENTAL SETTING

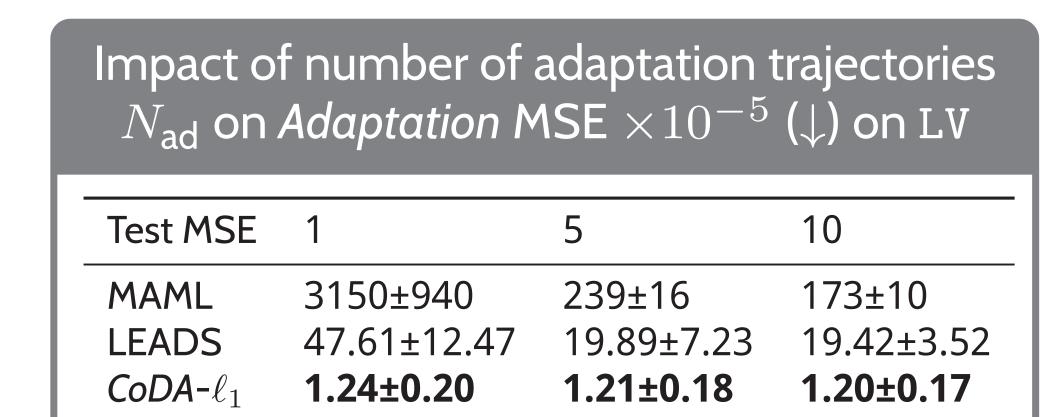
- ODE: Lotka-Volterra (LV), Glycolitic Oscillator (G0)
- PDE: Gray-Scott (GS), Navier-Stokes (NS)
- $d_p$  parameters vary between physical systems ( $d_p=2$  LV, GO, GS;  $d_p=1$ : NS)

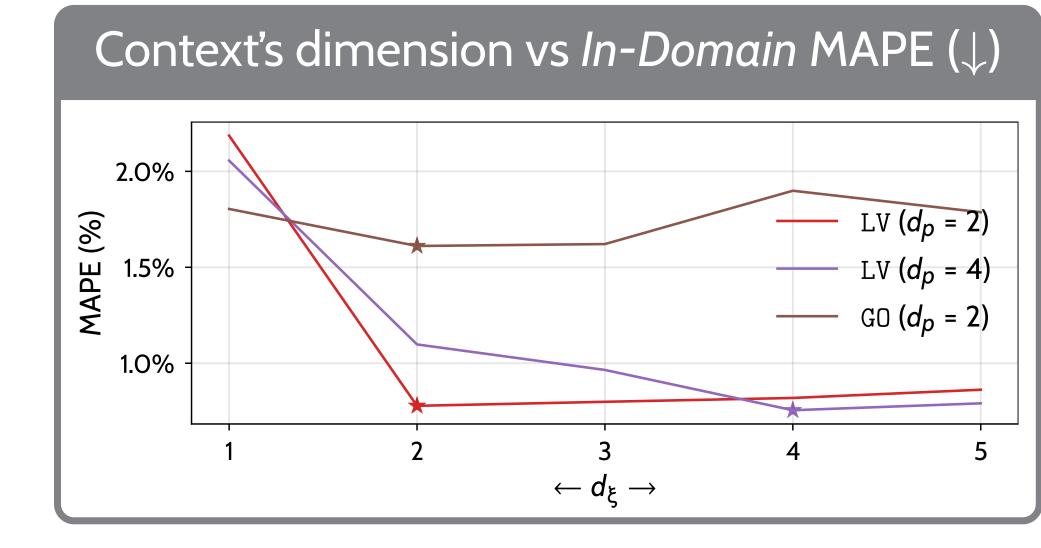
## RESULTS

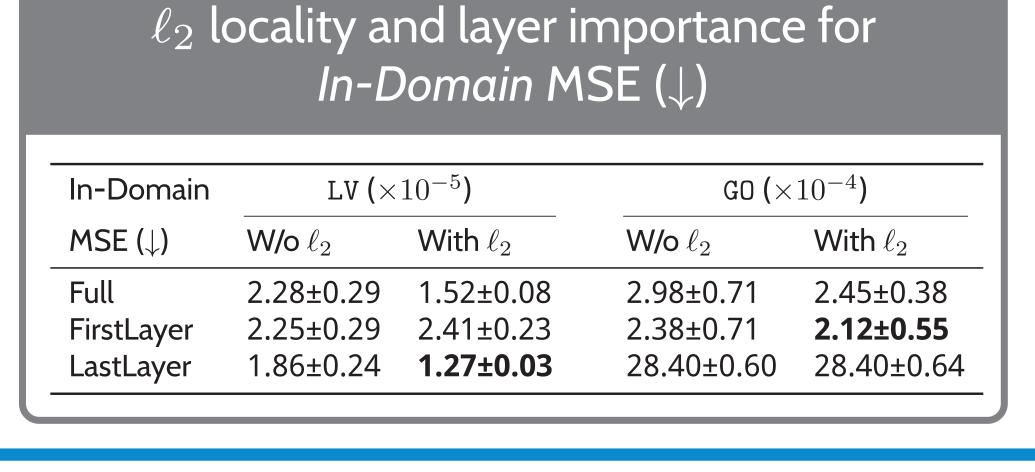












# SYSTEM PARAMETER ESTIMATION

- Learn correspondence o between context vectors oxplus and  $\mathit{known}$  system parameters oxplus on  $\mathcal{E}_{\mathrm{tr}}$
- Apply the correspondence to  $\mathcal{E}_{\mathrm{ad}}$  to infer *unknown* system parameters

