# Predictability Enables Parallelization of Nonlinear State Space Models
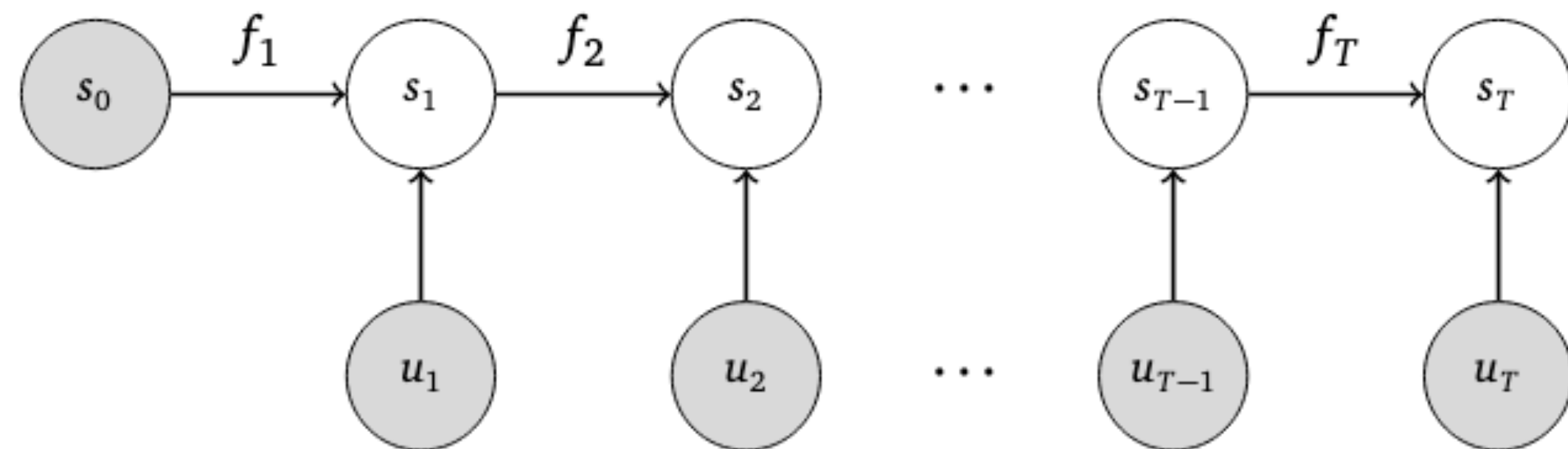
Xavier Gonzalez*, Leo Kozachkov*,

David M. Zoltowski, Kenneth L. Clarkson,
Scott W. Linderman

# Nonlinear State Space Models (nSSMs) are everywhere!

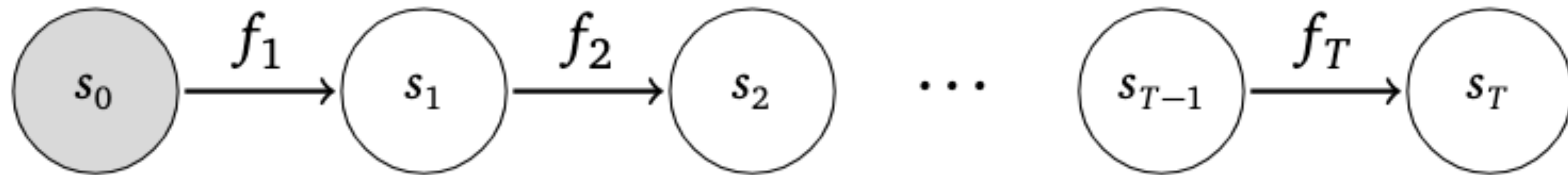$$s_t = f_t(s_{t-1}), \ \ s_t \in \mathbb{R}^D$$



We can include inputs by defining $f_t(s_{t-1}) = f(s_{t-1}, u_t)$



## Examples include

- Recurrent Neural Networks (RNNs)

- Markov chain Monte Carlo (MCMC)

- Sampling from a diffusion model

- The blocks of a transformer model *over depth*

- Gradient descent

- State in reinforcement learning agent

- …and more…

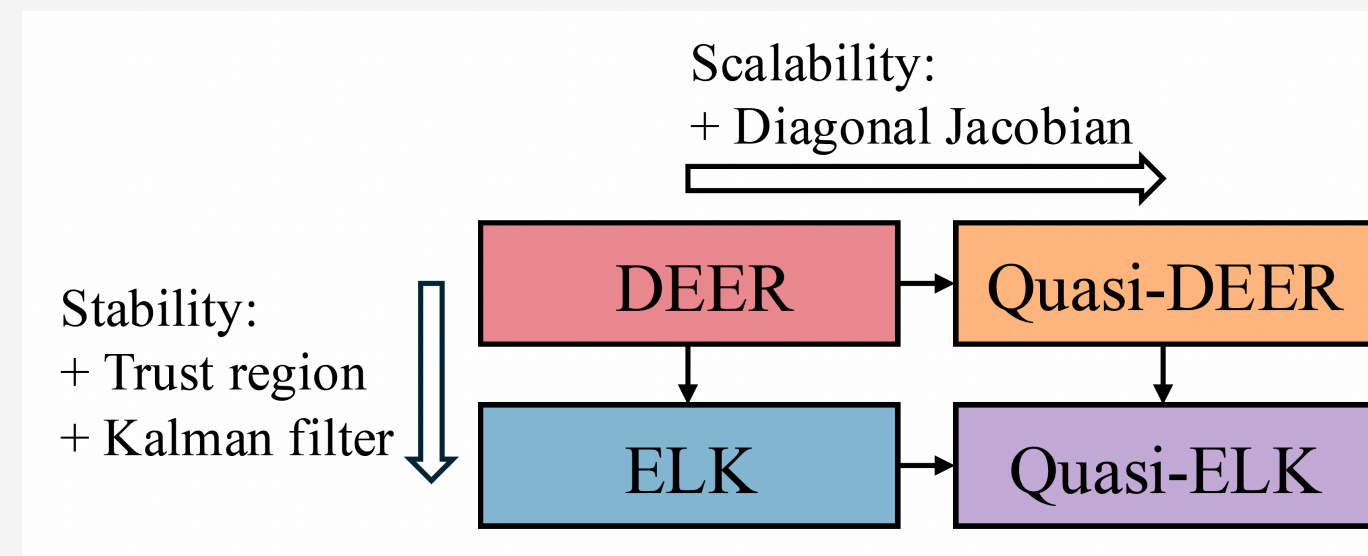# Evaluating nonlinear SSMs seemed to be "inherently sequential"…
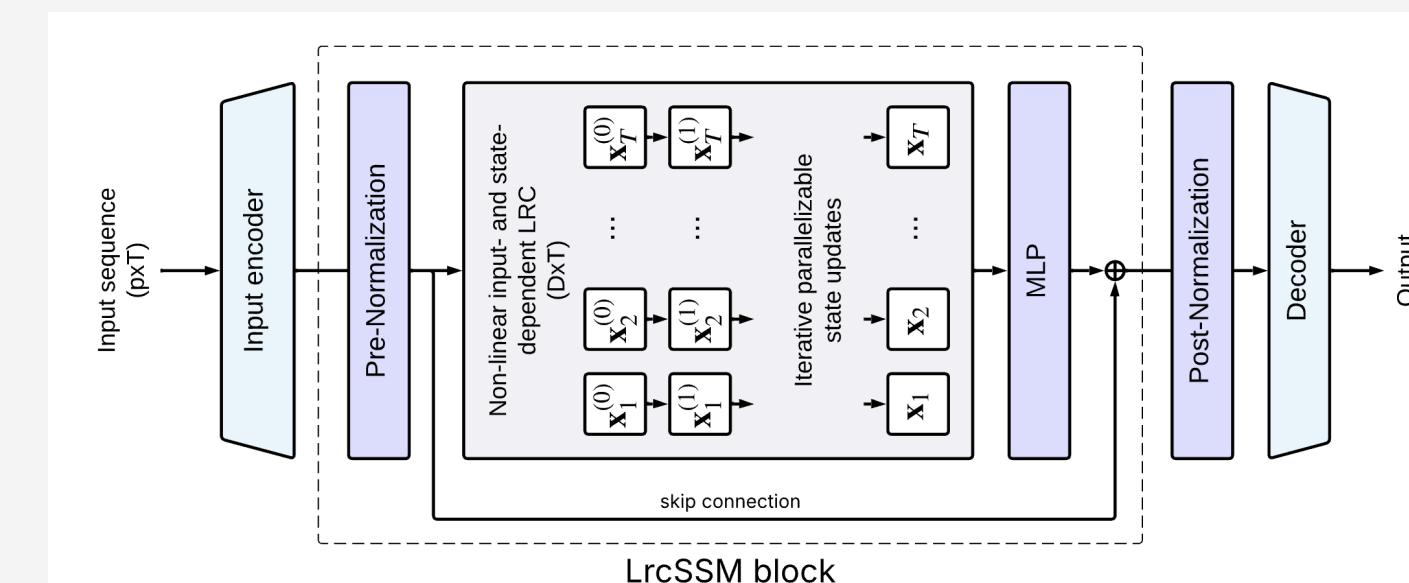
# …but turns out that nonlinear SSMs can be parallelized!
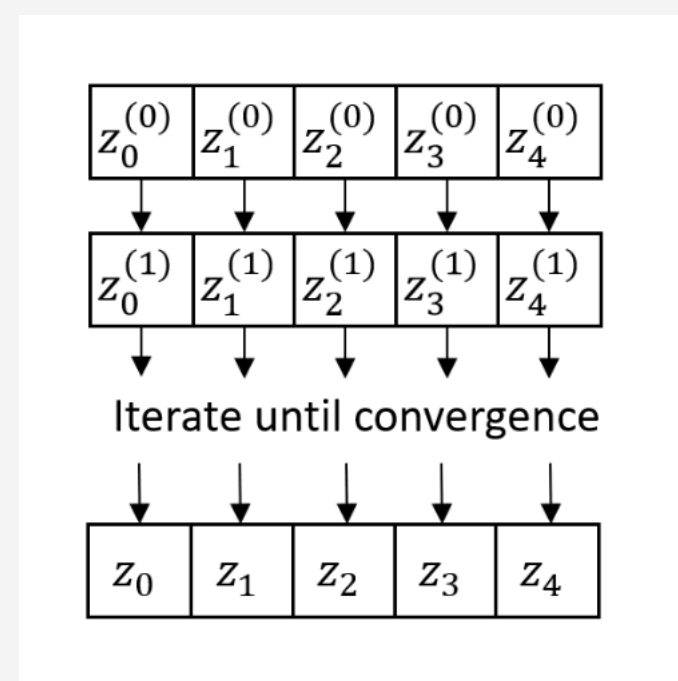
## Danieli et al, NeurIPS '23
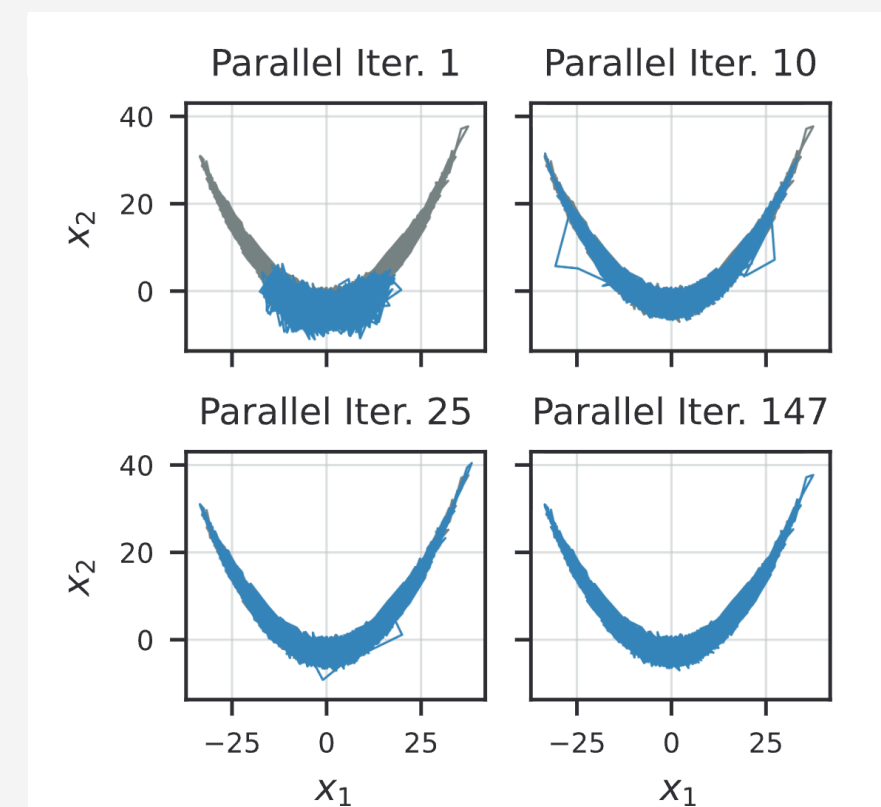


## Gonzalez et al, NeurIPS '24



## Farsang et al, NeurIPS '25
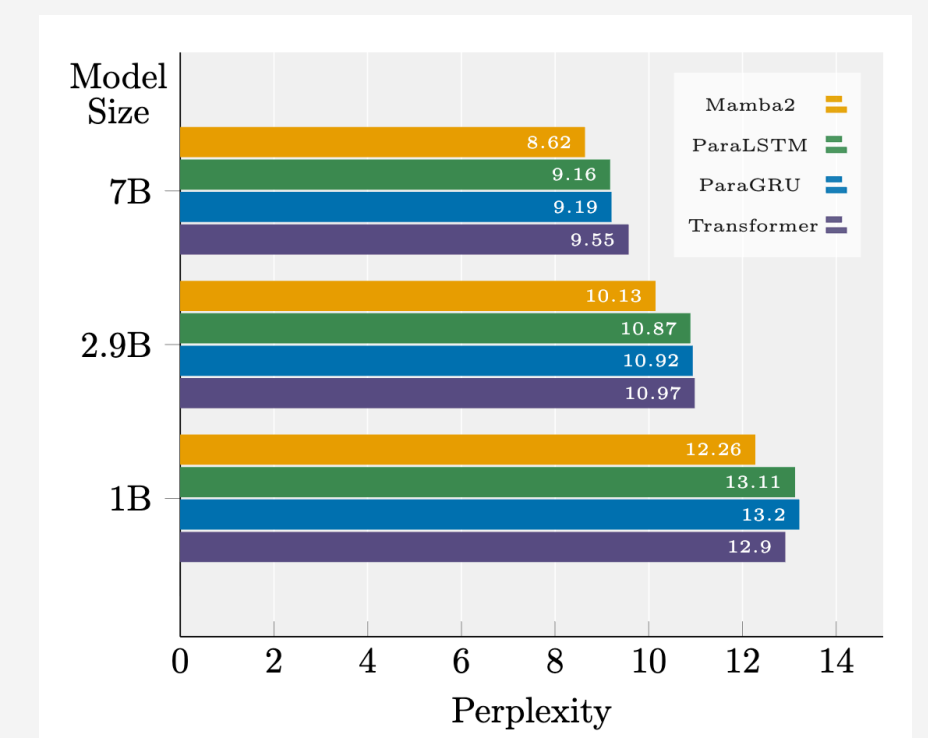


## Lim et al, ICLR '24



## Zoltowski et al, NeurIPS '25



## Danieli et al, '25

# We can parallelize nSSMs by going from <u>sequential evaluation</u> to a <u>high-dimensional optimization</u> problem

Recall the nSSM satisfies $s_t = f_t(s_{t-1})$, starting from some initial condition $s_0 \in \mathbb{R}^D$.

Treat as our optimization variable *the entire trajectory* $\mathbf{s}_{1:T} \in \mathbb{R}^{TD}$ .

Define the one-step prediction error as
$r_t := s_t - f_t(s_{t-1})$ .

Define the "merit function" as
$\mathscr{L}(\mathbf{s}_{1:T}) := \|r_1\|^2 + \|r_2\|^2 + \ldots + \|r_T\|^2.$

Then the unique minimizer of $\mathscr{L}$ is the true roll-out from the nSSM.

$\mathscr{L}(\mathbf{s}_{1:T})$

$s_{t-1}$

$\mathbf{s}_{1:T}^*$

$s_t$

# But how do we know when our solver will converge quickly?

We show that "predictability" enables parallelization.

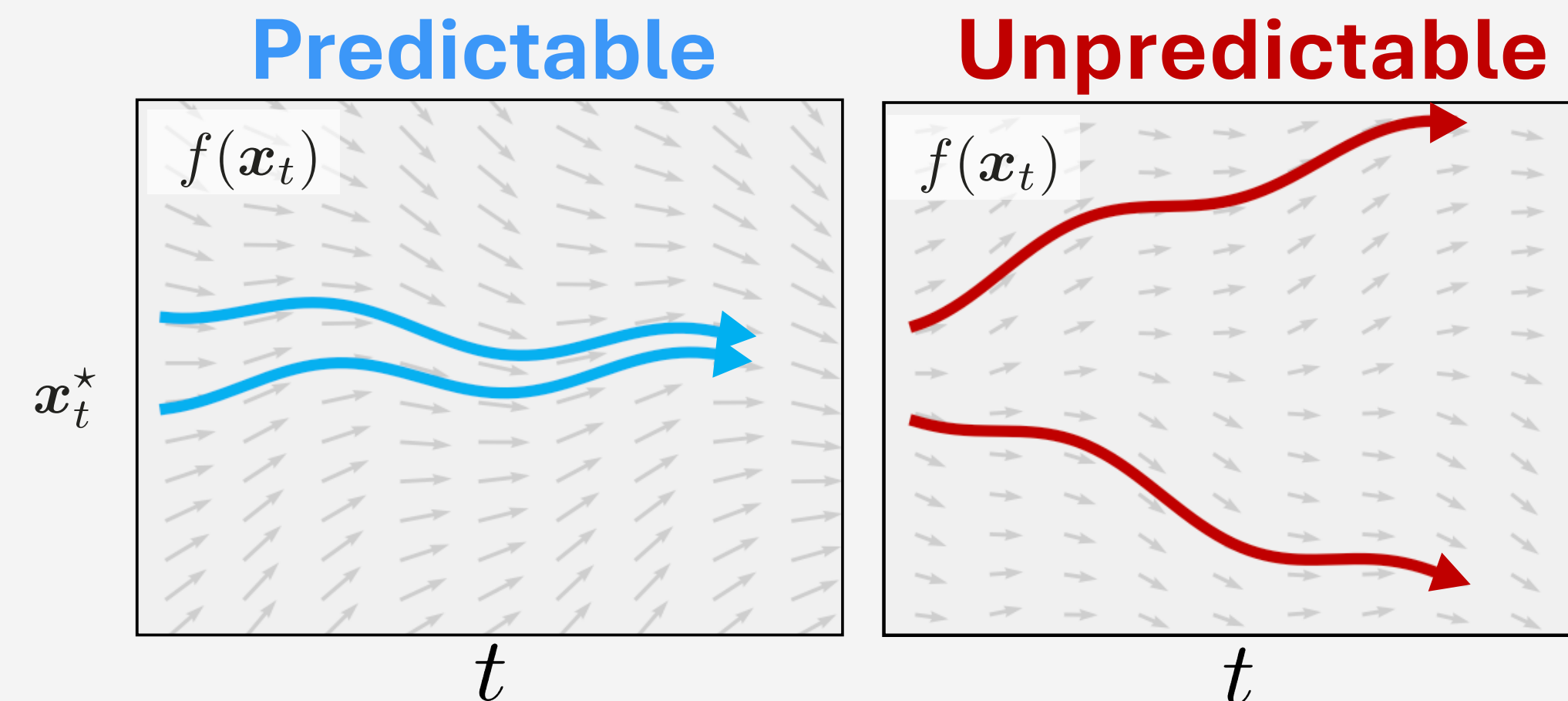We connect the "predictability" of the state-space *dynamics* to the conditioning of the *optimization problem.*



**Predictable**

**Unpredictable**

**Key quantity**

State-space Dynamics

$f(\boldsymbol{x}_t)$

$f(\boldsymbol{x}_t)$

$\boldsymbol{x}_t^\star$

$t$

$t$

Largest Lyapunov Exponent (LLE, or $\lambda$)

- $\lambda < 0$, predictable
- $\lambda > 0$, unpredictable

Optimization

$\mathcal{L}(\boldsymbol{x}_{1:T})$

$\mathcal{L}(\boldsymbol{x}_{1:T})$

$\boldsymbol{x}_t$

$\boldsymbol{x}_{t-1}$

$\boldsymbol{x}_{t-1}$

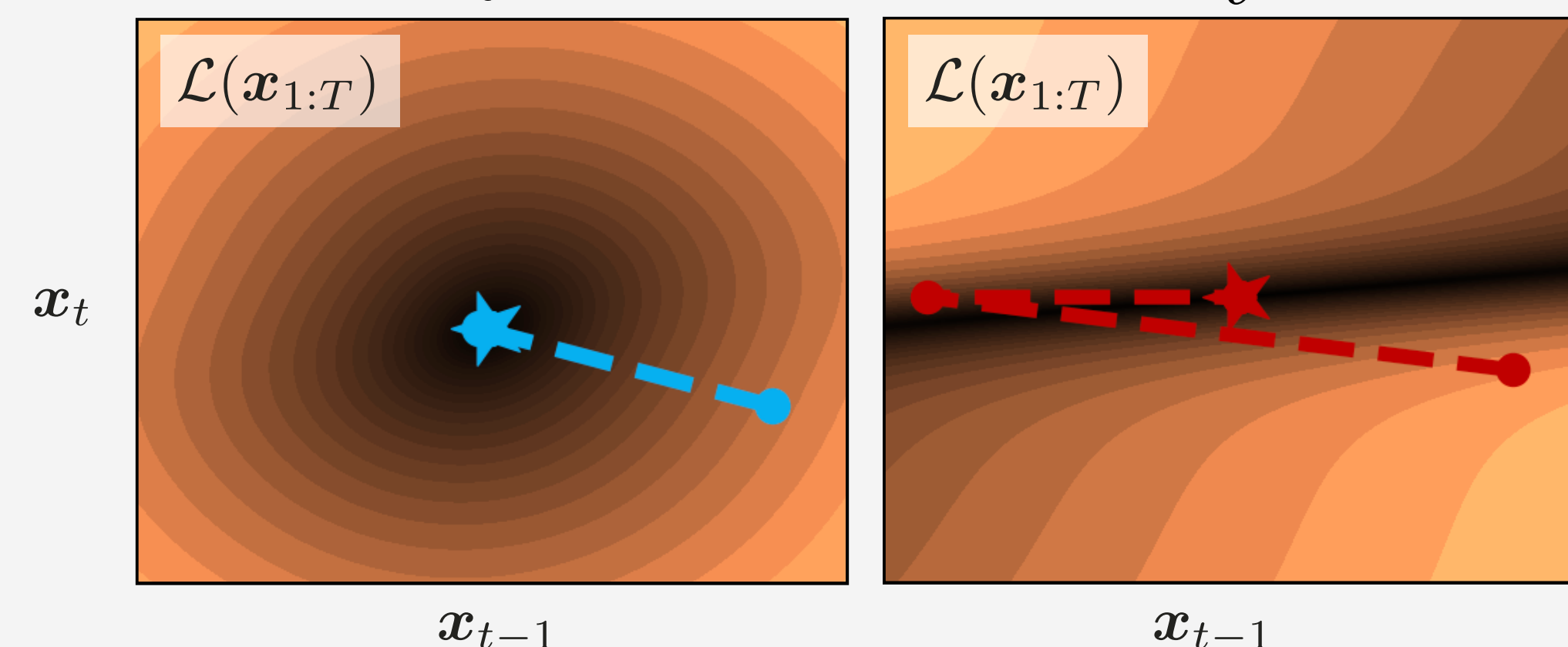PL constant $\mu$

- $\mu \gg 0$, well-conditioned, fast
- $\mu \approx 0$, ill-conditioned, slow

# Main Theorem

<u>Theorem 2 of our paper (informal):</u>

Let $\mu$ be the PL constant, representing the conditioning of the optimization problem. Larger $\mu$ is a better-conditioned optimization problem that can be solved more quickly.

Let $\lambda$ be the Largest Lyapunov Exponent (LLE), representing the predicability of the nSSM. $\lambda < 0$ means predicable dynamics, while $\lambda > 0$ means unpredictable dynamics.

Then, under regularity assumptions,

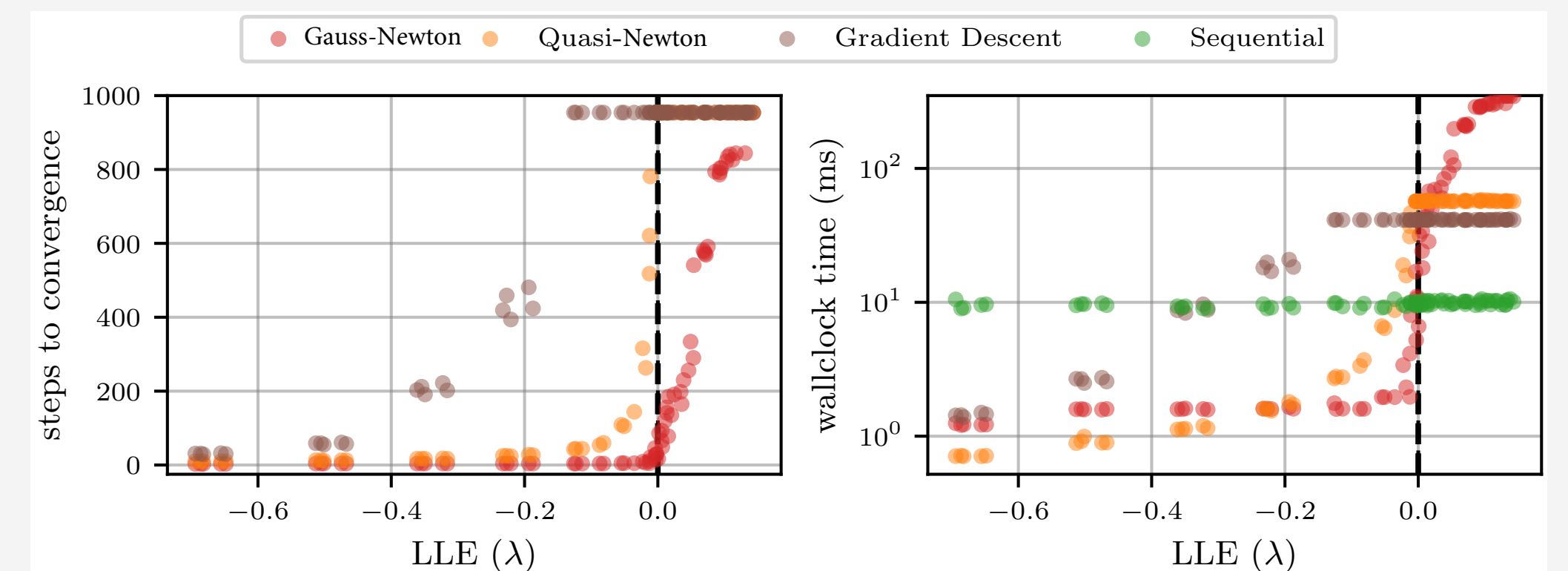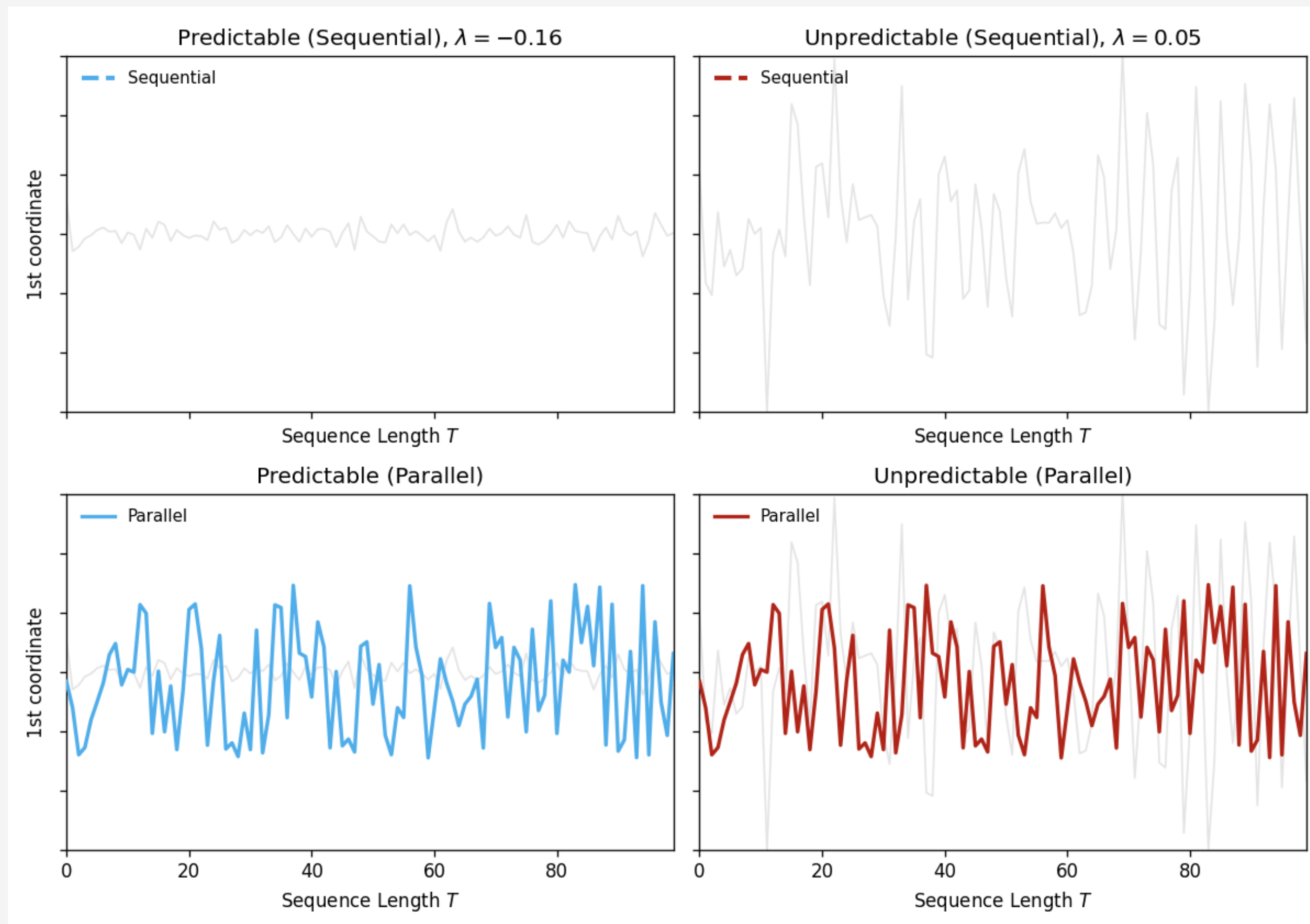$$\frac{e^\lambda - 1}{e^{\lambda T} - 1} \leq \sqrt{\mu} \leq \frac{1}{e^{\lambda(T-1)}} \; .$$

<span style="color:blue"><u>Predictable ($\lambda < 0$)</u></span>

- For long sequences, PL constant $\mu$ is bounded away from zero.
- Fast, well-conditioned optimization
- Speed-ups from parallelization

<span style="color:red"><u>Unpredictable ($\lambda > 0$)</u></span>

- For long sequences, PL constant $\mu \to 0$ goes to zero exponentially quickly
- Slow, ill-conditioned optimization
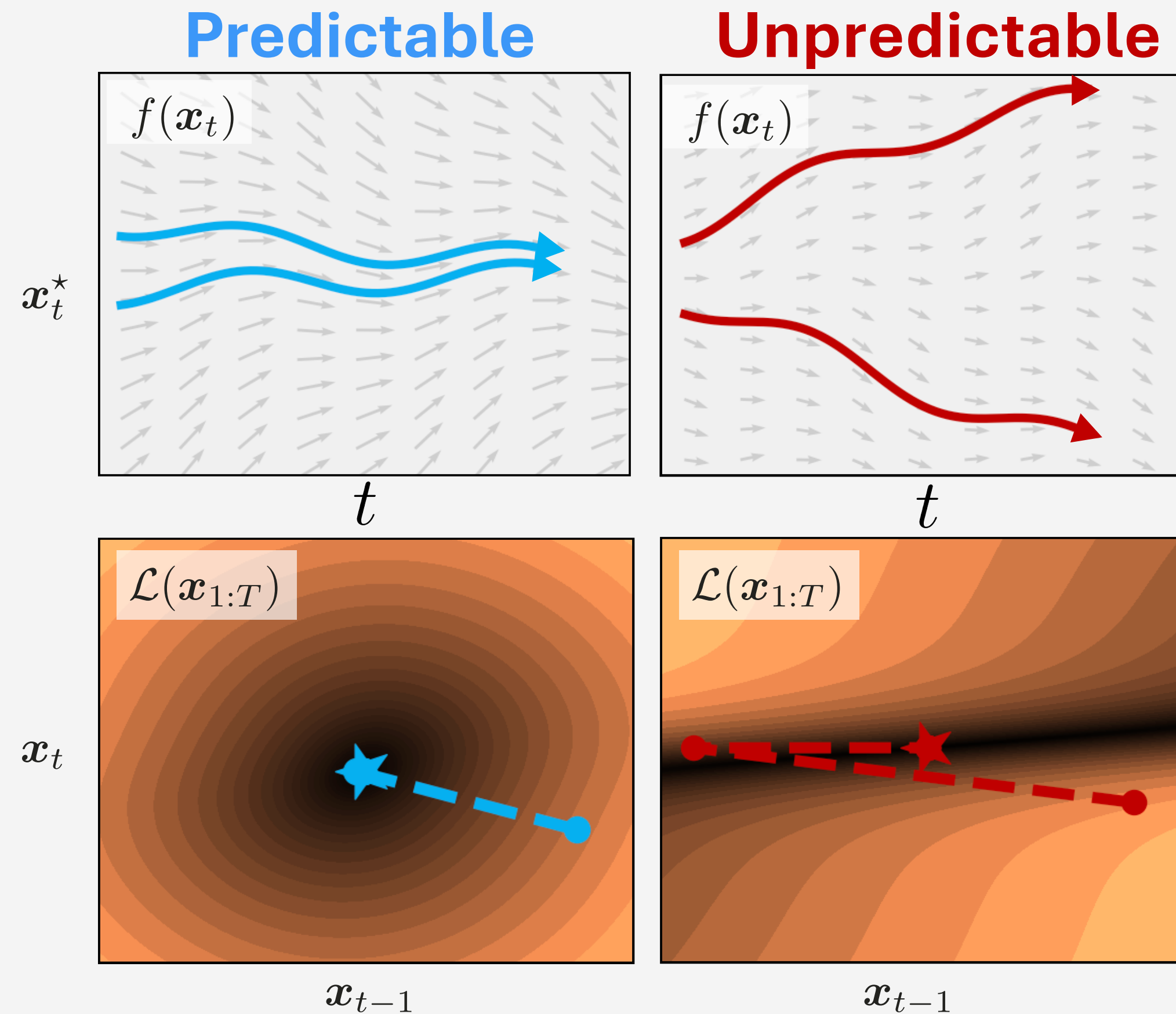- No speed-ups from parallelization :(

# Experiments

# Implications

- Predictable state space models are used in many applications (particularly in robotics and control), and <u>you should parallelize them!</u>

- If you want to build a sequence modeling architecture (RNN) so that it enjoys parallelized training, <u>you should design it to be predictable and therefore parallelizable.</u>

# Come by our poster to learn more!



- Paper: https://arxiv.org/abs/2508.16817

- Code: https://github.com/lindermanlab/predictability_enables_parallelization

- Longer Talk: https://www.youtube.com/watch?v=C9AqgW51-B4