

# The Three Regimes of Offline-to-Online Reinforcement Learning

**Lu Li, Tianwei Ni, Yihao Sun, Pierre-Luc Bacon**

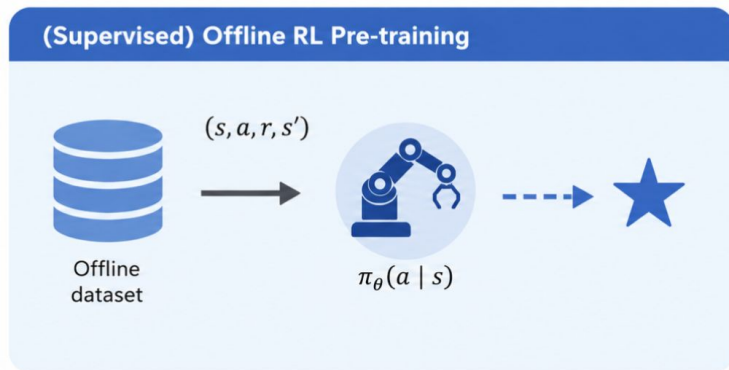
Lu Li at ICML 2026 Workshop on Decision-Making from Offline Datasets to Online Adaptation

July 11th, 2026



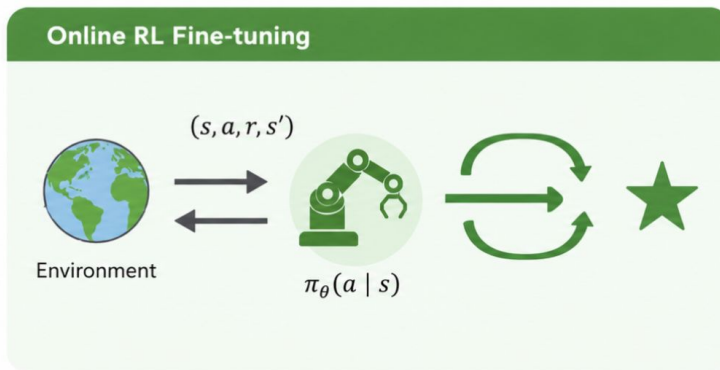
# Offline-to-Online RL

## Offline-to-Online RL



Offline pretraining:

$$\pi_0 = A_{\text{off}}(\mathcal{D})$$



Online fine-tuning:

$$\pi_N = A_{\text{on}}(\mathcal{M}, \mathcal{D}, \pi_0)$$

**Goal:** co-design offline algorithm and online algorithm,  
to maximize  $J(\pi_N)$

# Why Offline-to-Online RL?

- Online RL (from scratch) is highly **data-intensive**, making large-scale interaction impractical in many real-world applications.
- Offline RL leverages pre-collected datasets and avoids online interaction, but often produces **suboptimal policies**.

**Key idea of O2O RL:** Learn a policy from an offline dataset and subsequently fine-tune it through online interaction to further improve performance.

→ Combines the data efficiency of offline RL with the performance gains of online RL

# Inconsistency in Fine-tuning Results

# Prior Work

Warm-Start RL (WSRL): offline pretraining, only online data in fine-tuning

---

## Efficient Online Reinforcement Learning Fine-Tuning Need Not Retain Offline Data

Zhiyuan Zhou<sup>\*1</sup>, Andy Peng<sup>\*1</sup>, Qiyang Li<sup>1</sup>, Sergey Levine<sup>1</sup>, Aviral Kumar<sup>2</sup>

<sup>1</sup>UC Berkeley, <sup>2</sup>Carnegie Mellon University (\*Equal Contribution)

---

RL with Prior Data (RLPD): no offline pretraining, combining offline data and online data in fine-tuning

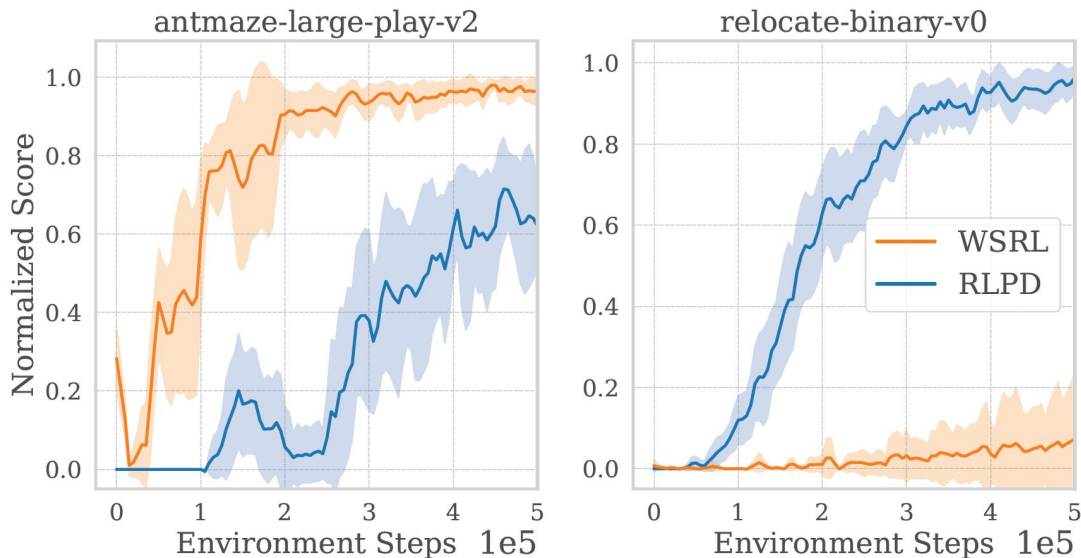
---

## Efficient Online Reinforcement Learning with Offline Data

---

Philip J. Ball<sup>\*1</sup> Laura Smith<sup>\*2</sup> Ilya Kostrikov<sup>\*2</sup> Sergey Levine<sup>2</sup>

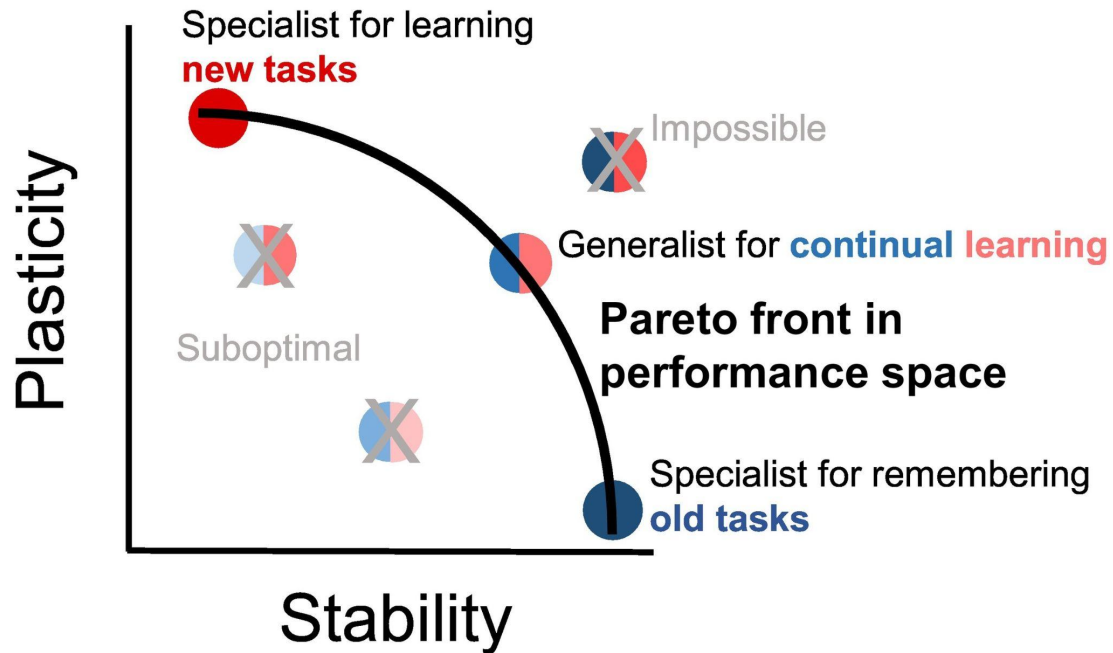
# Inconsistency in Online Fine-tuning



**Question:** What **underlying factors** cause design choices to succeed in some settings but fail in others?

# A Decomposition of Performance Based on Stability–Plasticity Principle

# Stability–Plasticity dilemma



Trends in Neurosciences

Figure from paper: metaplasticity to solving the Catastrophic Forgetting Problem

# Quantifying Priors

Performance of the pretrained policy:  $J(\pi_0)$

Performance of the the dataset:  $J(\pi_{\mathcal{D}}) \approx \frac{1}{N} \sum_{i=1}^N \sum_{t=1}^T r_{i,t}$ .

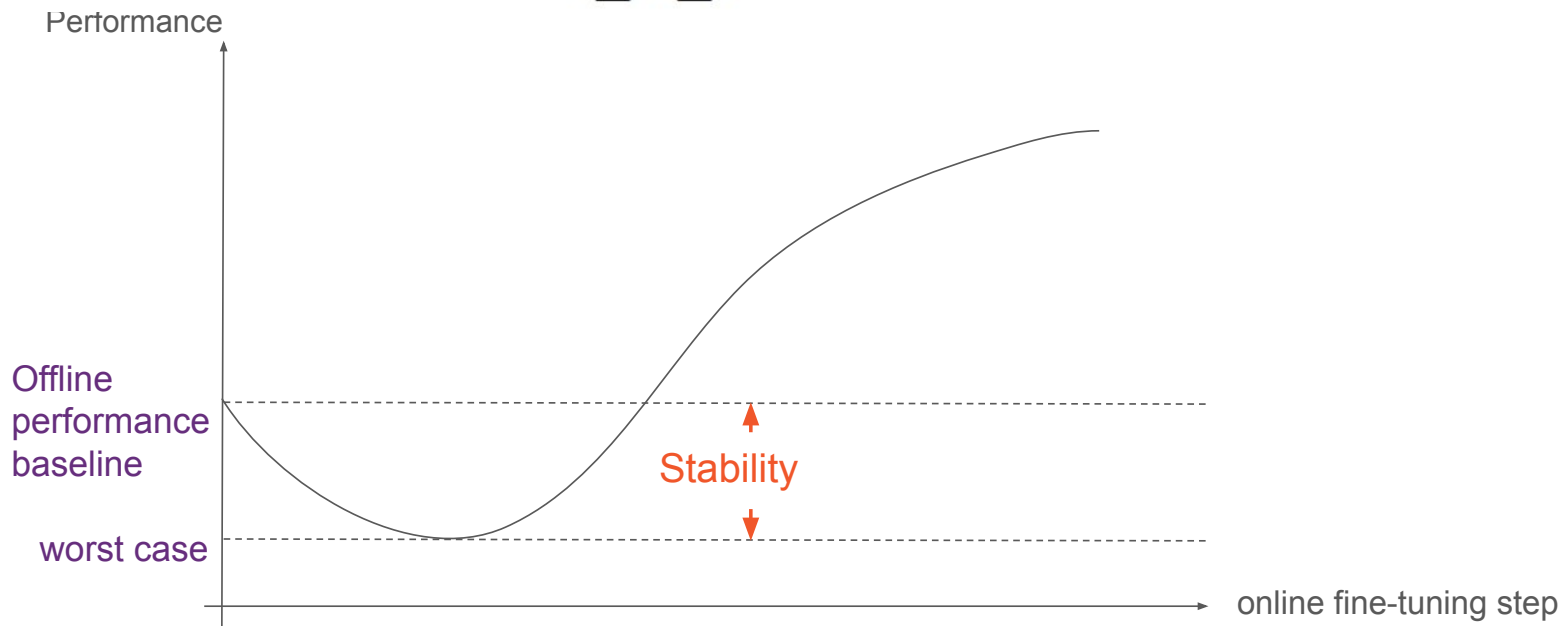
Offline performance baseline:

$$J_{\text{off}}^* := \max(J(\pi_0), J(\pi_{\mathcal{D}}))$$

# Performance Decomposition

Stability with respect to offline performance baseline:

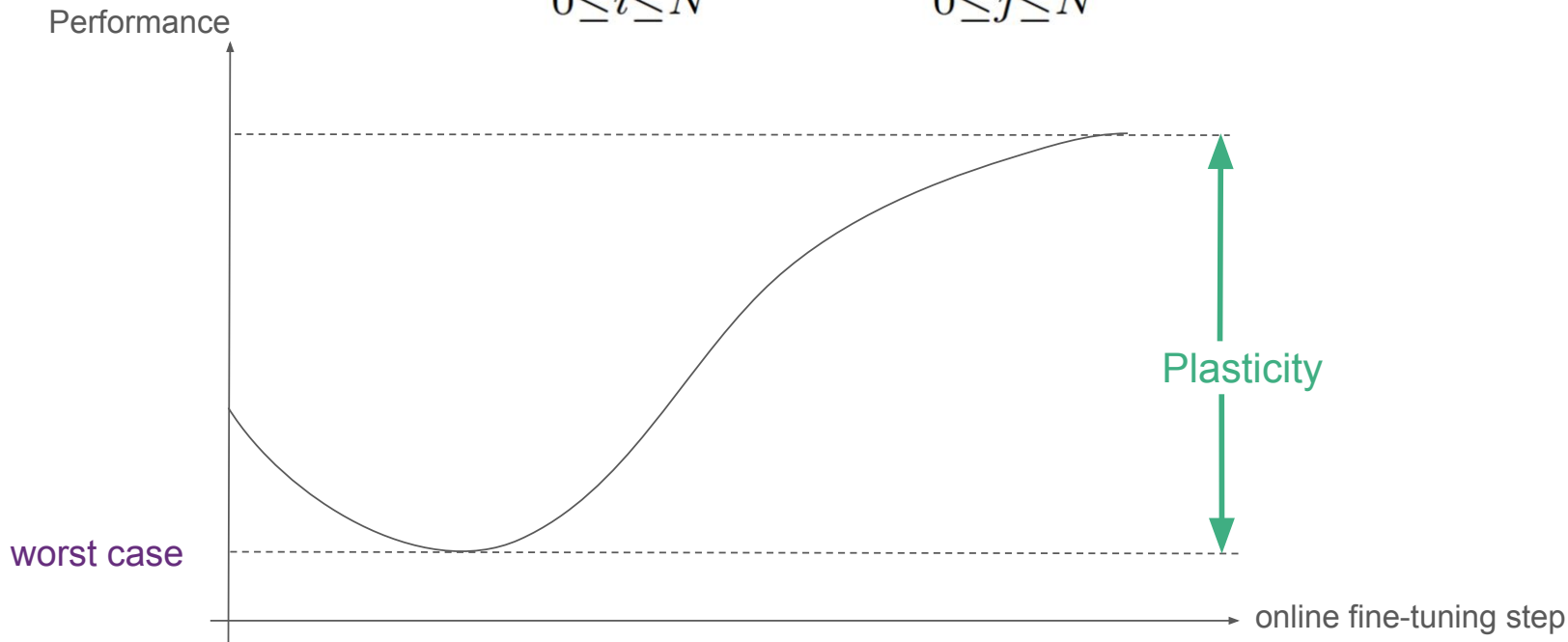
$$\text{Stability}(J_{\text{off}}^*) = \min_{0 \leq n \leq N} J(\pi_n) - J_{\text{off}}^* \leq 0.$$



# Performance Decomposition

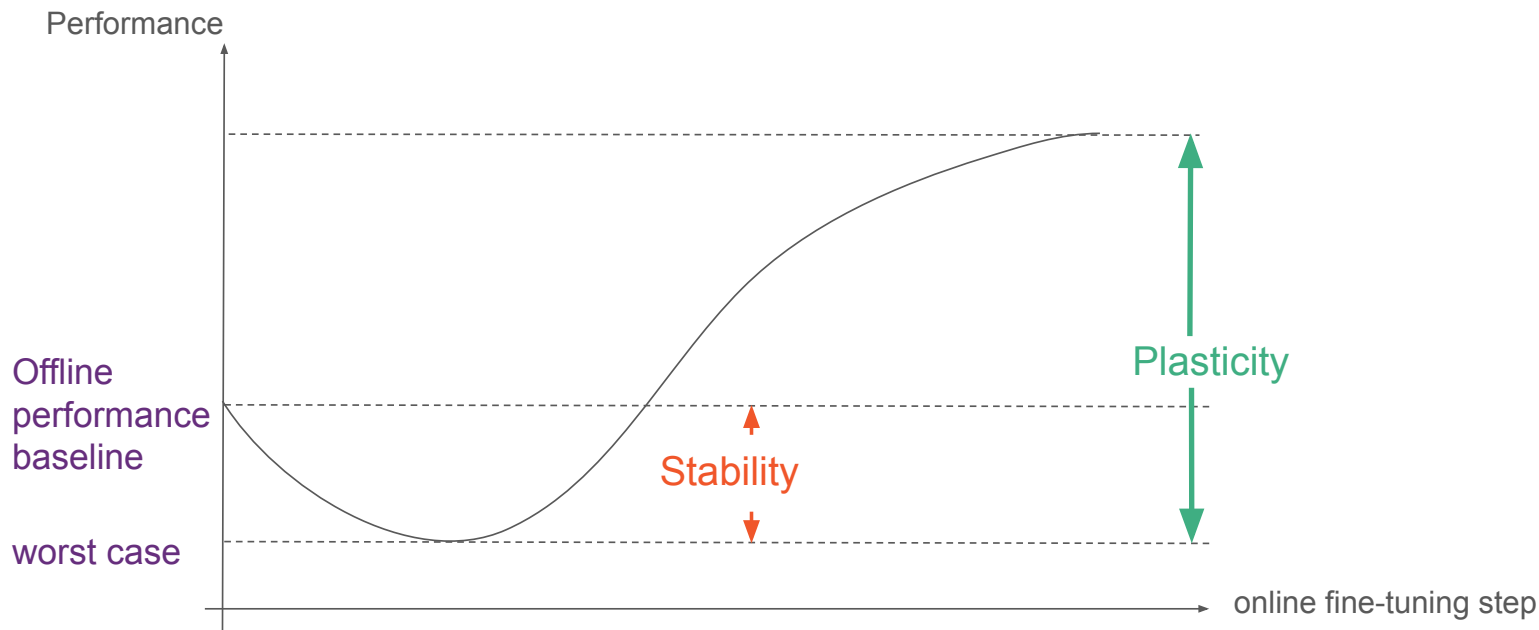
Plasticity: the ability to improve its performance, relative to the its lowest observed value

$$\text{Plasticity} = \max_{0 \leq i \leq N} J(\pi_i) - \min_{0 \leq j \leq N} J(\pi_j) \geq 0.$$

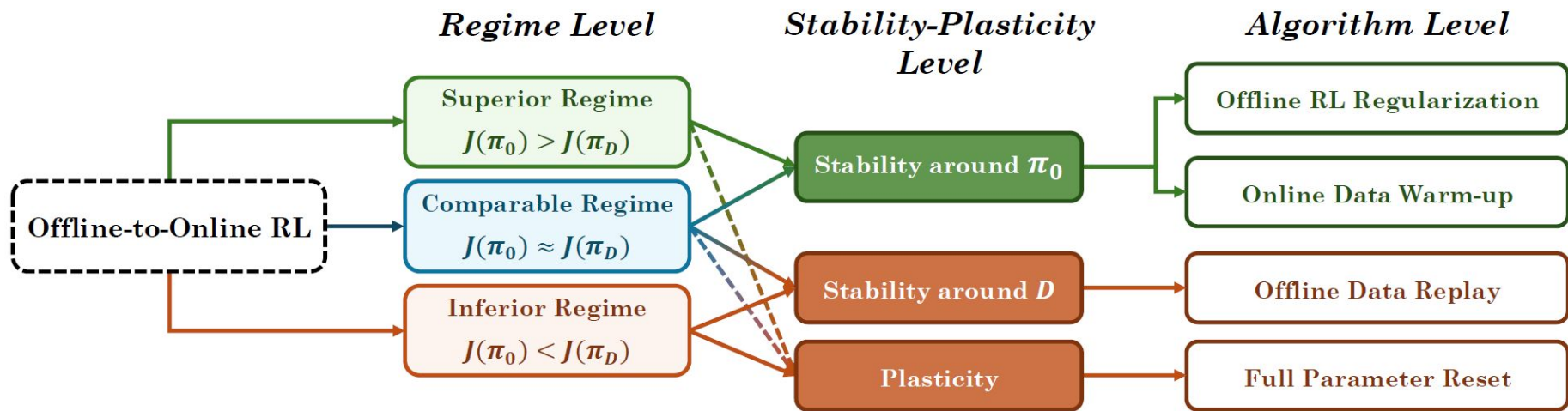


# Performance Decomposition

$$\underbrace{\max_{0 \leq n \leq N} J(\pi_n)}_{\text{Best Performance}} = \underbrace{J_{\text{off}}^*}_{(1) \text{ Prior}} + \underbrace{\text{Stability}(J_{\text{off}}^*)}_{(2) \text{ Degradation} \leq 0} + \underbrace{\text{Plasticity}}_{(3) \text{ Online Improvement} \geq 0}$$



# The three regimes



# Minimal baseline

We define a naive online RL fine-tuning baseline, which is ***intentionally minimalist***: use a standard online RL algorithm (e.g. SAC) initialized with an offline pretrained agent

- no offline data replay during fine-tuning
- no additional regularization
- ...

# Design Choices in Improving Stability

Stability around the pretrained policy

- online data warm-up
- offline RL regularization (e.g. conservatism)
- KL regularization
- learning rate warm-up
- ...

# Design Choices in Improving Stability

## Stability around the pretrained policy

- online data warm-up
- offline RL regularization (e.g. conservatism)
- KL regularization
- learning rate warm-up
- ...

## Stability around the offline dataset

- offline data replay(append): initial the buffer with offline data
- offline data replay(separate): one offline data buffer, one online buffer
- data generative methods
- ...

# Design Choices in Improving Plasticity

## Plasticity

- parameter reset: reset network parameters to random value, which severely degrades initial performance but also significantly enhances plasticity.
- other variant of parameter reset that proposed in online RL

Empirical Study

# Offline Pretraining

Use 21 task-dataset compositions from D4RL dataset,  
with 3 different pretraining algorithms:

- CalQL: offline RL algorithm based on SAC
- ReBRAC: offline RL algorithm based on TD3
- Behavior Cloning: imitation learning algorithm (commonly used in robotics)

lead to 63 settings

# Regime Classification

use T-test with margin to classify the pretrained results into three regimes: **Superior**, **Inferior**, and **Comparable**

$H_0 : \mu_0 - \mu_D \leq \delta$  vs  $H_1 : \mu_0 - \mu_D > \delta$     If rejected  $\rightarrow$  **Superior**

$H_0 : \mu_D - \mu_0 \leq \delta$  vs  $H_1 : \mu_D - \mu_0 > \delta$     If rejected  $\rightarrow$  **Inferior**

If neither test rejects  $\rightarrow$  **Comparable**

# Categorizing Algorithms into 4 Types

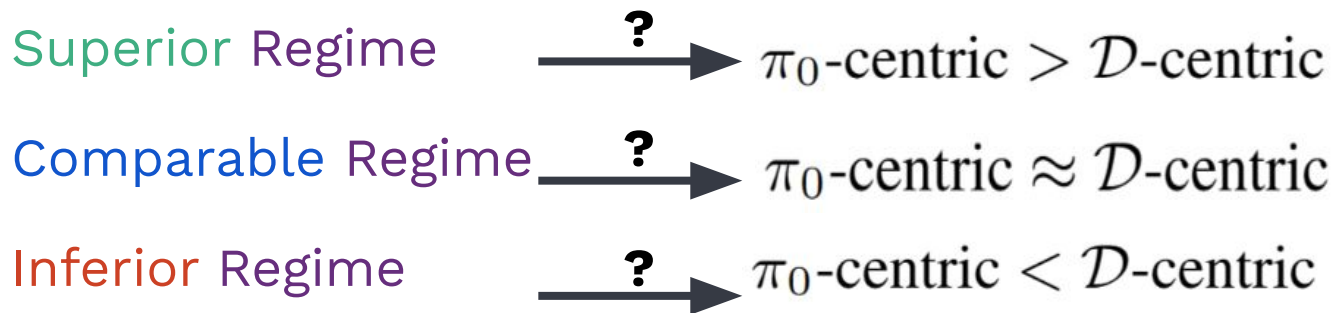
group algorithms by the primary source of stability they emphasize:

- minimal baseline
  - standard online RL (SAC or TD3)
- $\pi_0$ -centric method
  - online data warm-up (WSRL-style)
  - offline RL regularization (CalQL or ReBRAC)
- $D$ -centric method
  - offline data replay (use online and offline buffer with 1:1 ratio)
  - offline data replay + full parameter reset (RLPD-style)
- mixed  $\pi_0+D$  method
  - offline RL regularization + offline data replay (Most prior O2O work)

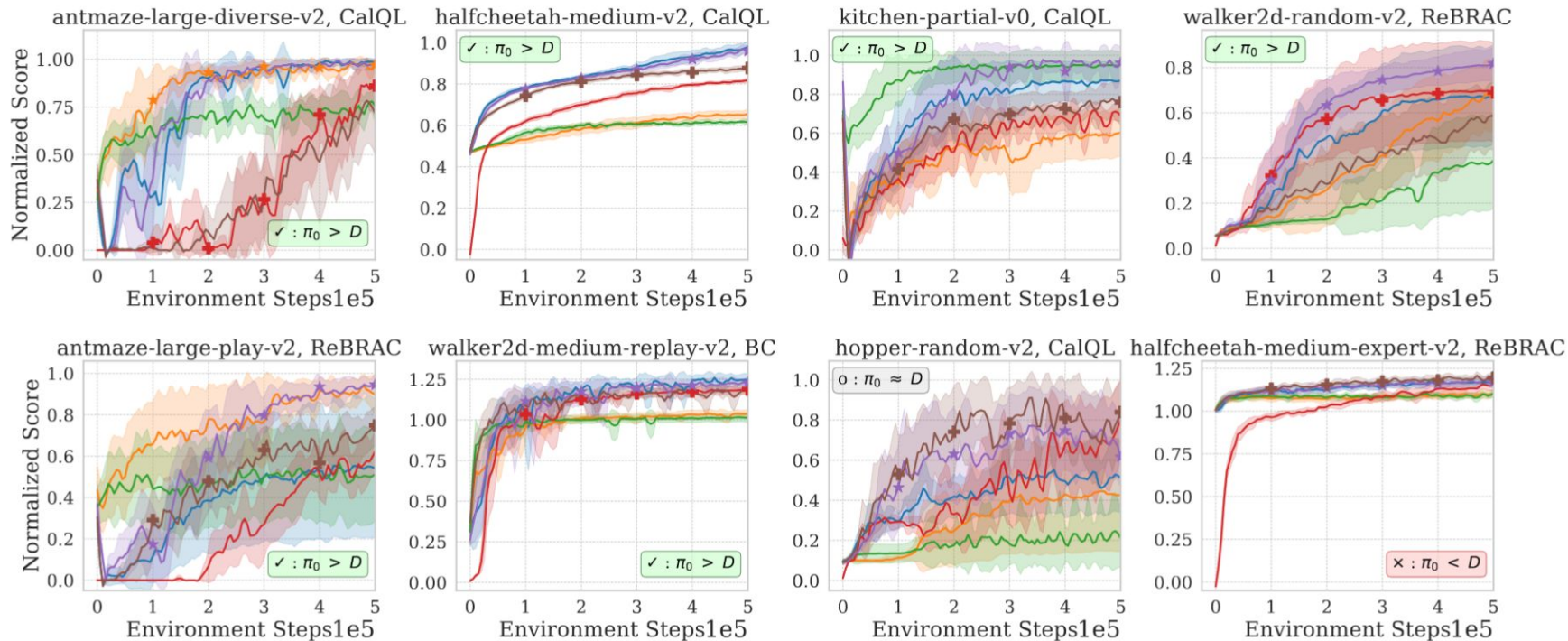
# Comparing Fine-tuning Results

We focus on compare the strongest  $\pi_0$ -centric methods and  $\mathcal{D}$ -centric methods to better approximate the ideal performance achievable by each stability source.

Thus, there are three kinds of fine-tuning results, which is obtained by t-test:



# Superior Regime



# Fine-tuning Results in Superior Regime

		Pretraining Regime		
		Superior	Comparable	Inferior
Fine-tune	$\pi_0$ -centric $>$ $\mathcal{D}$ -centric	24	2	1
	$\pi_0$ -centric $\approx$ $\mathcal{D}$ -centric	6	2	3
	$\pi_0$ -centric $<$ $\mathcal{D}$ -centric	2	4	19

24/32(**75%**) correct predictions, only 2/32(**6%**) opposite mismatches

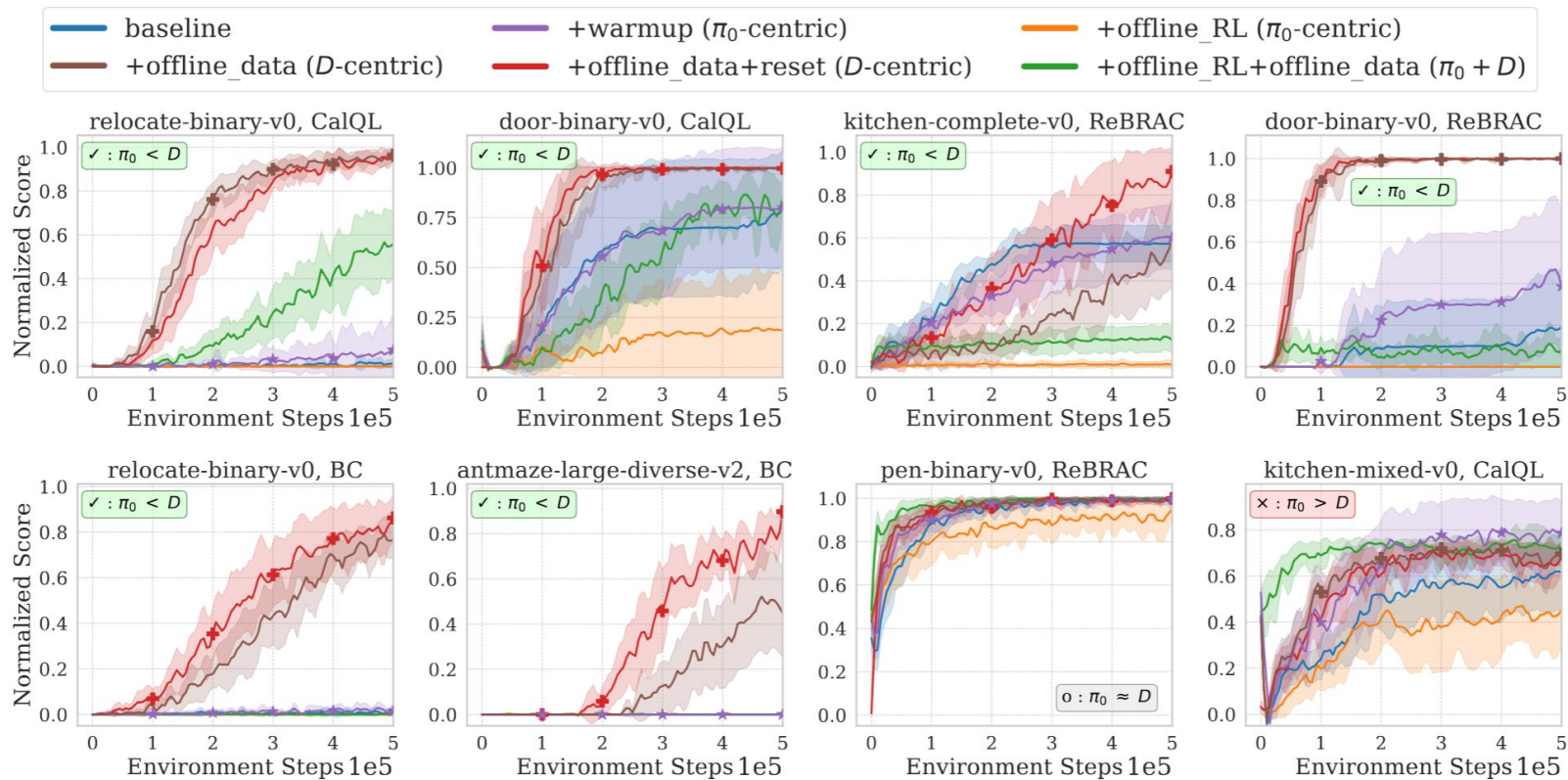
**Takeaway:** In the Superior regime,  $\pi_0$ -centric methods typically provide more effective fine-tuning than  $\mathcal{D}$ -centric methods.

# Fine-tuning Results in Superior Regime

Within  $\pi_0$ -centric methods,  
*online data warm-up* outperform *offline RL regularization* in 27/32(84%).  
*offline RL regularization* (and *+offline data replay*) only works well in those settings where  $\pi_0$  is close to optimal.

Since *offline RL regularization* provide much stronger stability which at the cost of plasticity, lead to less performance degradation, but also limits further improvements.

# Inferior Regime



# Fine-tuning Results in **Inferior** Regime

		Pretraining Regime		
		Superior	Comparable	Inferior
Fine-tune	$\pi_0$ -centric $>$ $\mathcal{D}$ -centric	24	2	1
	$\pi_0$ -centric $\approx$ $\mathcal{D}$ -centric	6	2	3
	$\pi_0$ -centric $<$ $\mathcal{D}$ -centric	2	4	19

19/23(**83%**) correct predictions, only 1/23(**4%**) opposite mismatches

**Takeaway:** In the **Inferior** regime,  $\mathcal{D}$ -centric methods typically provide more effective fine-tuning than  $\pi_0$ -centric methods.

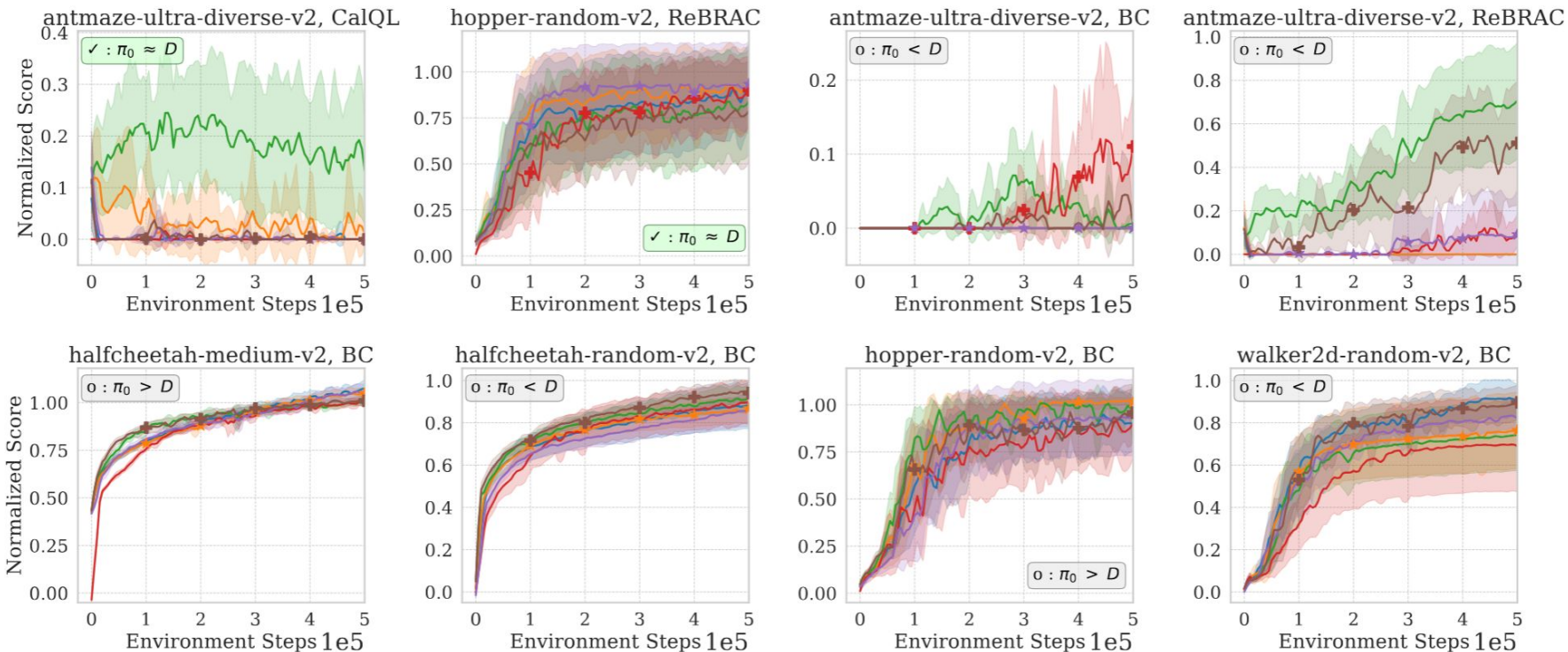
# Fine-tuning Results in **Inferior** Regime

Within  $D$ -centric methods,

Out of 24 settings, *offline data replay + reset* performs better in 13 settings and *offline data replay* performs better in 11 settings.

resetting the parameters allows the agent to adapt and acquire new knowledge more effectively when offline pretraining phase substantially reduces plasticity while offering limited knowledge

# Comparable Regime



# Fine-tuning Results in Comparable Regime

		Pretraining Regime		
		Superior	Comparable	Inferior
Fine-tune	$\pi_0$ -centric $>$ $\mathcal{D}$ -centric	24	2	1
	$\pi_0$ -centric $\approx$ $\mathcal{D}$ -centric	6	2	3
	$\pi_0$ -centric $<$ $\mathcal{D}$ -centric	2	4	19

2/8(25%) correct predictions.

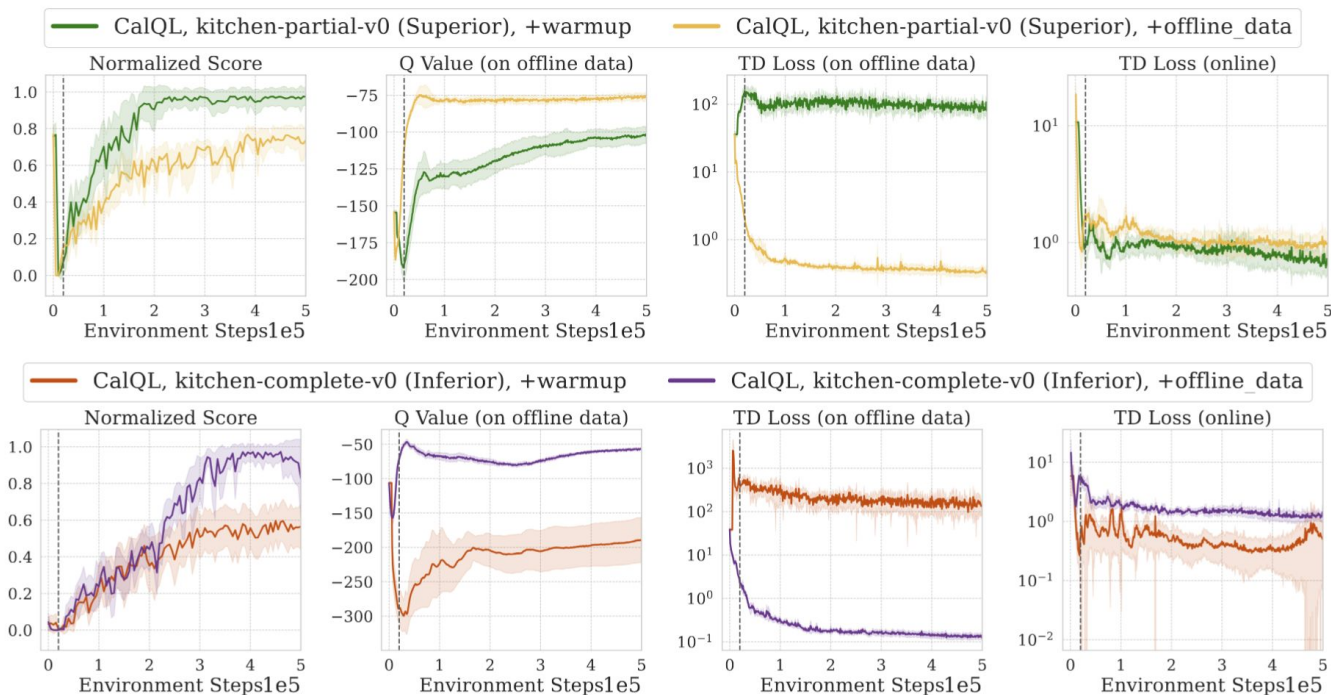
**Takeaway:** In the **Comparable** regime,  $\pi_0$ -centric method and  $\mathcal{D}$ -centric method should in principle yield comparable fine-tuning outcomes, though in practice their relative performance is often sensitive to implementation details.

# Overall Results and Alternative Metrics

		Pretraining Regime		
		Superior	Comparable	Inferior
Fine-tune	$\pi_0$ -centric $>$ $\mathcal{D}$ -centric	24	2	1
	$\pi_0$ -centric $\approx$ $\mathcal{D}$ -centric	6	2	3
	$\pi_0$ -centric $<$ $\mathcal{D}$ -centric	2	4	19

Domains	Metric	Accuracy $\uparrow$	Opposite mismatch $\downarrow$
Sparse-reward (27 cases)	Raw return (ours)	<b>78%</b>	<b>4%</b>
	Dense-reward proxy	65%	11%
All domains (63 cases)	Raw return (ours)	<b>71%</b>	<b>5%</b>
	Q-function	51%	30%
	BC performance	41%	<b>5%</b>

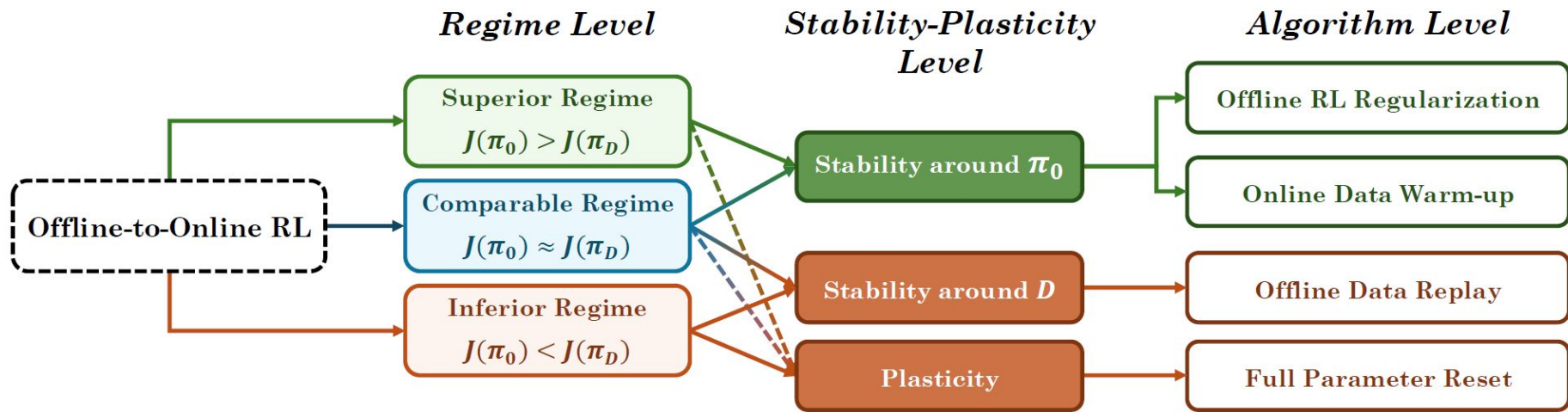
# Mechanistic Analysis



Fine-tuning without offline data results in significantly higher TD loss on the offline dataset in the **Inferior** regime compared to the **Superior** regime, and also leads to divergence of the Q-values on the offline data.

Conclusion

# The three regimes



- We provide a **clear explanation** of the conflicting empirical evidence: design choices that seem inconsistent across benchmarks in fact reflect different underlying regimes.
- We offer an **actionable guidance** for practitioners. By identifying the regime of a given setting, one can select methods that align with its stability–plasticity requirements, reducing reliance on trial-and-error.
- The three regimes provides an efficient lens for **understanding offline-to-online RL**. But we also have to recognize that such discretizing behavior is a simplification.

# Thanks to co-authors



**Tianwei Ni**

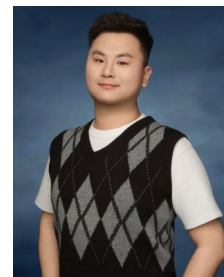


**Yihao Sun**



**Pierre-Luc Bacon**

**Special thanks to Guozheng Ma for insightful discussions.**



Thanks for listening!  
Looking forward to your thoughts!

contact:  
lu.li@mila.quebec  
twini2016@gmail.com



Read our paper



Read our blog