

# Zero-shot Model-based Reinforcement Learning using Large Language Models

Abdelhakim Benechehab<sup>12</sup>, Youssef Attia El Hili<sup>1</sup>, Ambroise Odonnat<sup>13</sup>,  
Oussama Zekri<sup>4</sup>, Albert Thomas<sup>1</sup>, Giuseppe Paolo<sup>1</sup>, Maurizio Filippone<sup>5</sup>,  
Ievgen Redko<sup>1</sup>, Balázs Kégl<sup>1</sup>

<sup>1</sup> Huawei Noah's Ark Lab, Paris, France

<sup>2</sup> Department of Data Science, EURECOM

<sup>3</sup> Inria, Univ. Rennes 2, CNRS, IRISA

<sup>4</sup> ENS Paris-Saclay

<sup>5</sup> Statistics Program, KAUST

—

March 30, 2025



**1** Preliminaries

2 Problem setup

3 Approach

4 Results

5 Conclusion



Reinforcement Learning environments are Markov decision processes  $\mathcal{M} = \langle \mathcal{S}, \mathcal{A}, P, r, \mu_0, \gamma \rangle$ , where:

- $\mathcal{S}$  state space,  $\mathcal{A}$  action space.
- Transition fn  $P_t : (s, a, s') \mapsto \mathbf{Pr}(s_{t+1} = s' | s_t = s, a_t = a)$ .
- Reward function  $r : (s, a) \mapsto r(s, a)$ .
- $\mu_0$  initial state distribution,  $\gamma \in [0, 1]$  discount factor.





The goal of RL is to find a policy  $\pi : \mathcal{S} \rightarrow \Delta(\mathcal{A})$  that maximizes the return:

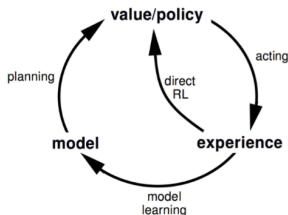
$$\eta(\pi) := \mathbb{E}_{s_0 \sim \mu_0, a_t \sim \pi, s_{t>0} \sim P_t} \left[ \sum_{t=0}^{\infty} \gamma^t r(s_t, a_t) \right]$$



**Model-based RL (MBRL)** learns the transition  $\hat{P}$  from interaction data. The model maximizes the log-likelihood:

$$\mathcal{L}(\mathcal{D}; \hat{P}) = \frac{1}{N} \sum_{i=1}^N \log \hat{P}(s_{t+1}^i | s_t^i, a_t^i)$$

The learned model is used for policy search under the *learned* MDP  $\hat{\mathcal{M}} = \langle \mathcal{S}, \mathcal{A}, \hat{P}, r, \mu_0, \gamma \rangle$ .



towards data science - <https://tinyurl.com/3kxkx4p4>



Large Language Models (LLMs) are **transformer**-based, **decoder only** models trained using **autoregressive** next token prediction.

## LLaMA 3 Tokenizer

- Digits: ['0', '1', ... '999']
- Token Ids: [15, 16, 17, ... 5500]



"151,167,...,267"

## Time series processing

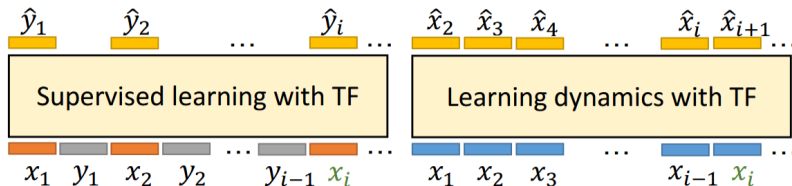
- Time series: [0.2513, 5.2387, 9.7889]
- Rescale+Encode: [150, 516, 850]
- Input str: '150,516,850,'
- Input str token list: ['150', ',', '516', ',', '850', ',', '']
- Input str token Id list: [3965, 11, 20571, 11, 16217, 11]

LLMTime

Sampling: *Softmax* over the digits tokens



In-context learning	Input prompt	Desired Output
Natural language processing	berry, baya, apple, manzana, <b>banana</b>	plátano
	Japan, mochi, France, croissant, <b>Greece</b>	baklava
Supervised learning $y_i = f(x_i) + \text{noise}$	$x_1, y_1, x_2, \dots, x_{i-1}, y_{i-1}, \mathbf{x_i}$	$f(x_i)$
Dynamical systems $x_{i+1} = f(x_i) + \text{noise}$	$x_1, x_2, x_3, \dots, x_{i-2}, x_{i-1}, \mathbf{x_i}$	$f(x_i)$



Transformers as Algorithms: Generalization and Stability in In-context Learning



1 Preliminaries

2 Problem setup

3 Approach

4 Results

5 Conclusion





- State space  $\mathbb{R}^{d_s}$ , Action space  $\mathbb{R}^{d_a}$ , Reward  $\mathbb{R}$
- Given a trajectory

$$\tau^\pi = (s_0, a_0, s_1, a_1, s_2, \dots, s_{T-1})$$

We want to learn the distribution of the next state using ICL and a pre-trained LLM with parameters  $\theta$ :

$$\{\hat{P}_\theta^{\pi,j}(s_t^j | \tau^\pi)\}_{t \leq T, j \leq d_s} = \text{ICL}_\theta(\tau^\pi)$$

Challenges:

- 1 Multivariate states:  $d_s > 1$
- 2 Including **actions** in-context:  $P(s_t^j | s_0, a_0, s_1, a_1, s_2, \dots, s_{T-1})$



1 Preliminaries

2 Problem setup

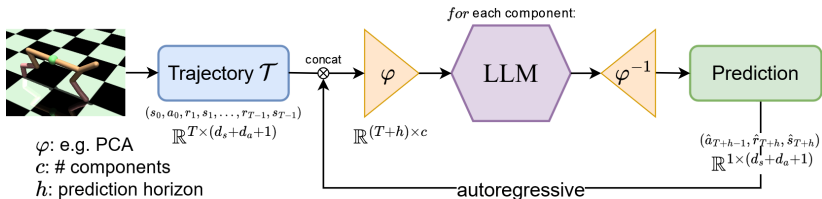
**3 Approach**

4 Results

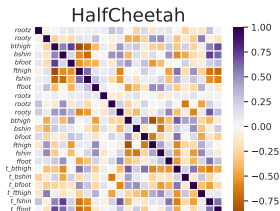
5 Conclusion



## DICL: Disentangled In-Context Learning [1]



In practice, we project states and actions  $(s, a)$  into the space of PCA components.





1 Preliminaries

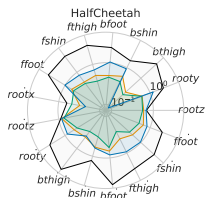
2 Problem setup

3 Approach

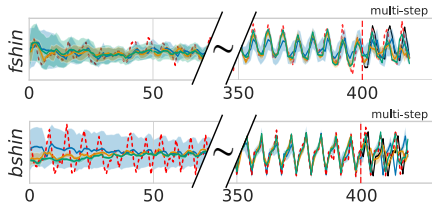
**4 Results**

5 Conclusion

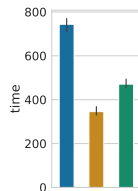
# Results: Prediction Error



Multi-step error



Predicted trajectories



Time

..... groundtruth    — vICL    — DICL-(s)    — DICL-(s, a)    — MLP

**PCA-based DICL achieves smaller multi-step error in less computational time.** We compare **DICL-(s)** and **DICL-(s, a)** using a number of components equal to half the number of features, with the vanilla approach **vICL** and an MLP baseline.



SAC: Soft Actor-Critic (an off-shelf RL algorithm)

+

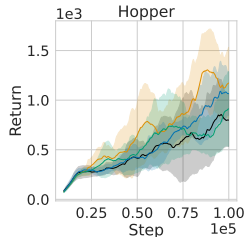
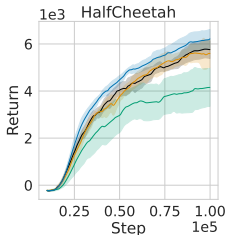
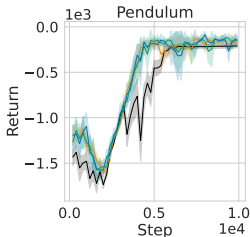
DICL

=

**DICL-SAC**

```

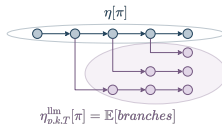
for  $t = 1, \dots, N\_interactions$  do
  New transition  $(s_t, a_t, r_t, s_{t+1})$  from  $\pi_\theta$ 
  Add  $(s_t, a_t, r_t, s_{t+1})$  to  $\mathcal{R}$ 
  Store auxiliary action  $\tilde{a}_t \sim \pi_\theta(\cdot|s_t)$ 
  if Generate LLM data then
    Sample trajectory  $\mathcal{T} = (s_0, \dots, s_{T_{\max}})$  from  $\mathcal{R}$ 
     $\{\hat{s}_{i+1}\}_{0 \leq i \leq T_{\max}} \sim \text{DICL-}(s)(\mathcal{T})$ 
    Add  $\{(s_i, \hat{a}_i, r_i, \hat{s}_{i+1})\}_{T \leq i \leq T_{\max}}$  to  $\mathcal{R}_{llm}$ 
  end if
  if update SAC then
    Sample batch  $\mathcal{B}$  of size  $b$  from  $\mathcal{R}$ 
    Sample batch  $\mathcal{B}_{llm}$  of size  $\alpha \cdot b$  from  $\mathcal{R}_{llm}$ 
    Update  $\phi$  and  $\psi$  on  $\mathcal{B} \cup \mathcal{B}_{llm}$ 
  end if
end for
    
```



— SAC — Dicl-SAC( $\alpha = 5\%$ ) — Dicl-SAC( $\alpha = 10\%$ ) — Dicl-SAC( $\alpha = 25\%$ )



Under mild assumptions on the LLM prediction error  $\varepsilon_{llm}$ , we have:



## Theorem (Multi-branch return bound)

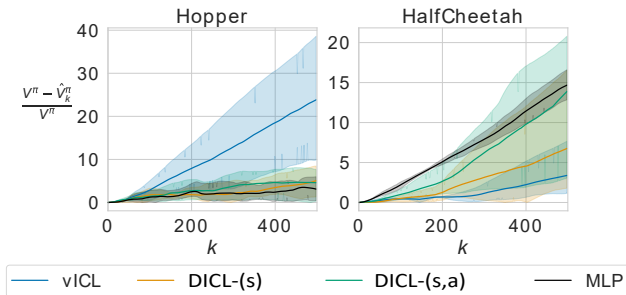
- $T$  the context length
- $p \in [0, 1]$  probability of branching
- $k$  the branch length
- $\varepsilon_{llm}$  the LLM in-context learning prediction error

$$|\eta(\pi) - \eta_{p,k,T}^{llm}(\pi)| \leq 2 \frac{\gamma^T}{1 - \gamma} r_{\max} k^2 p \varepsilon_{llm}(T)$$

where  $r_{\max} = \max_{s \in \mathcal{S}, a \in \mathcal{A}} r(s, a)$ .



We leverage **DICL** for **Policy evaluation** by predicting the rewards trajectory in-context.

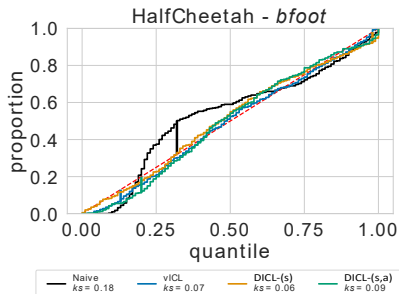






**Quantile calibration:** For probabilistic regression, a perfectly calibrated forecaster means that  $p\%$  of groundtruth values fall within the  $p\%$ -confidence interval of the predicted CDF.

LLMs are well-calibrated  
in-context forecasters.





1 Preliminaries

2 Problem setup

3 Approach

4 Results

5 Conclusion




- We presented **DICL**, a methodology to adapt LLMs for the task of dynamics learning in MBRL.
- We leveraged **DICL** for data-augmented RL, policy evaluation, and showed that LLMs are well-calibrated.

## Take Home Message

**LLMs** are powerful foundation models trained on vast amounts of data

→ **DICL** is an effective way to adapt them to MBRL



-  A. Benechehab, Y. A. E. Hili, A. Odonnat, O. Zekri, A. Thomas, G. Paolo, M. Filippone, I. Redko, and B. Kégl, “Zero-shot model-based reinforcement learning using large language models,” in *The Thirteenth International Conference on Learning Representations (ICLR)*, 2025.

# Thank You!

Want to know more?



Contact:

<https://abenechehab.github.io/>

✉ [abdelhakim.benechehab@gmail.com](mailto:abdelhakim.benechehab@gmail.com)

🔄 🐦 **in** @abenechehab

Slides available at:

<https://abenechehab.github.io/assets/pdf/dicl.pdf>