

# Data-Centric Knowledge-Enhanced Reasoning and Alignment of Large Language Models

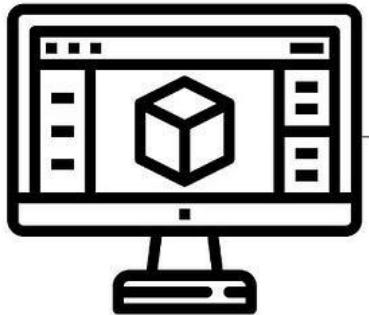
Xiusi Chen  
UIUC-NLP Talk  
Oct. 25, 2024

# About Me

- Xiusi (Hugh) Chen
- Postdoc @ Blender Lab
- Before UIUC: Ph.D. in CS @ UCLA
  - Thesis Title: One Step towards Autonomous AI Agents: Reasoning, Alignment and Planning
  - Thesis Committee: Wei Wang, Yizhou Sun, Kai-Wei Chang, Jeff Brantingham

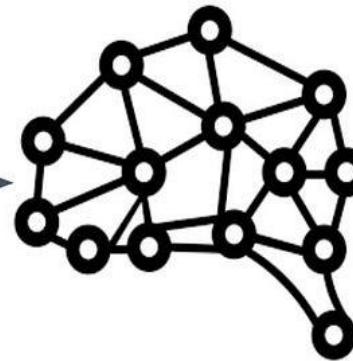
# How does AI Benefit Society?

*Where We Were*



Existing Software

*Where We Are*



LLMs and Diffusion Models  
(Foundation Models)

*Where We're Going*



AI Agents

# Core Properties of AI Agents

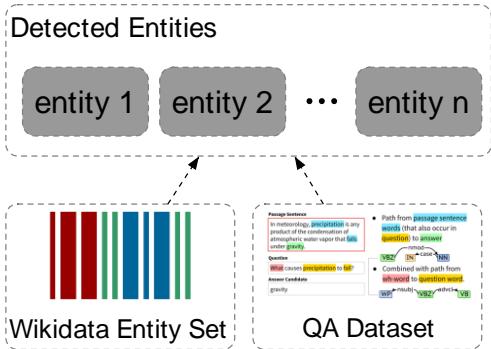
Strong Reasoning Ability

Well Aligned to Human  
Preference and Values

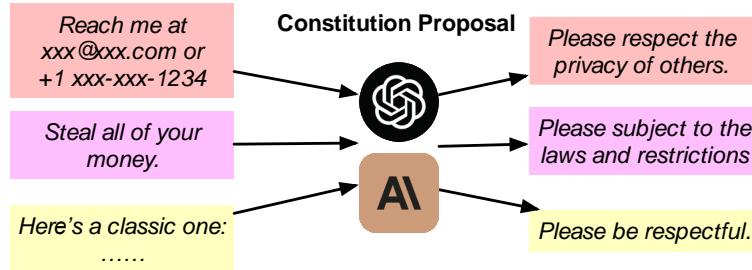
Planning Ahead

# My Research: Overview

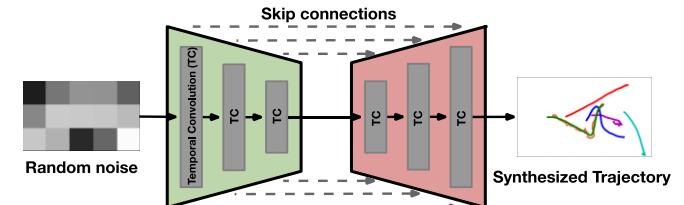
## Part I: Data-Centric Knowledge-Enhanced Reasoning



## Part II: Automatic Constitution Discovery and Self-Alignment

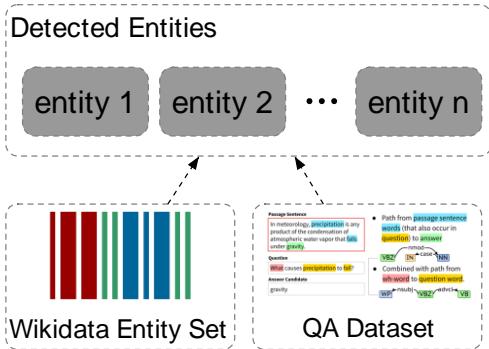


## Part III: Dynamics Modeling and Agents Planning

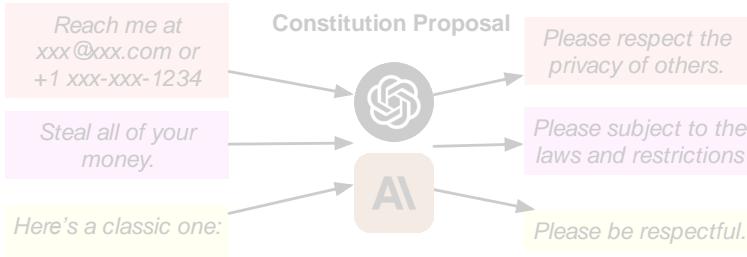


# My Research: Part I

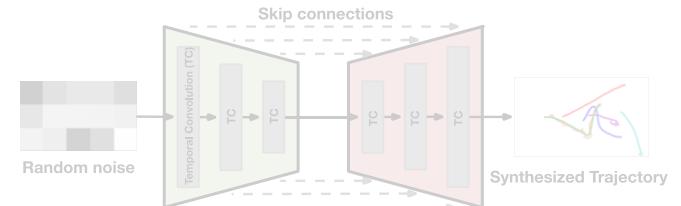
## Part I: Data-Centric Knowledge-Enhanced Reasoning



## Part II: Automatic Constitution Discovery and Self-Alignment



## Part III: Dynamics Modeling and Agents Planning



# Limitations of Pre-trained Language Models (PLMs)

Factual Error

“Albert Einstein won the Nobel Prize in Chemistry”

Logical Error

“If you add two apples to two oranges, you get four oranges.”

Bias and Discrimination

Generating text that implies certain ethnicities are inherently less intelligent or more prone to criminal behavior.

Privacy Violations

“XXX’s home address is \*\*\*, phone number is \*\*\*”

# Minimally-Supervised Data Generation and Selection

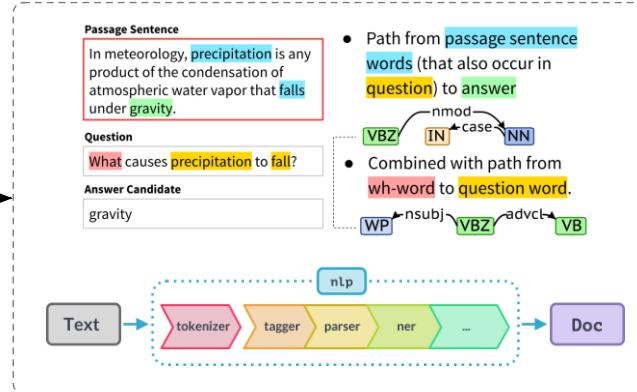
- Pre-training
  - Language and knowledge understanding
  - Costly, massive raw text
  - Most people use pre-trained LMs
- Fine-Tuning
  - Task adaptation
  - Smaller and focuses on a particular domain or task
  - Efficiency matters to broader users

# Data-Centric Knowledge-Enhanced Reasoning

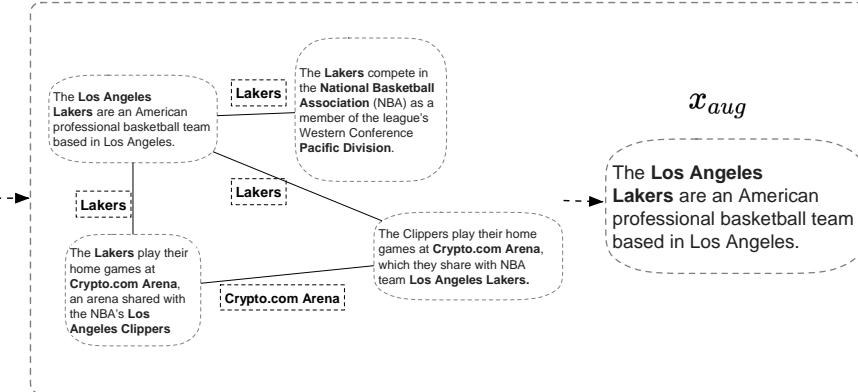
## QA data Acquisition



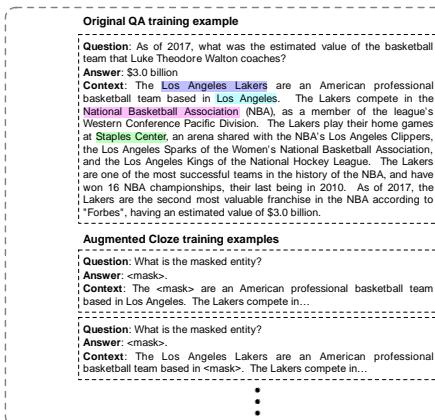
## Named Entity Recognition & Entity Typing



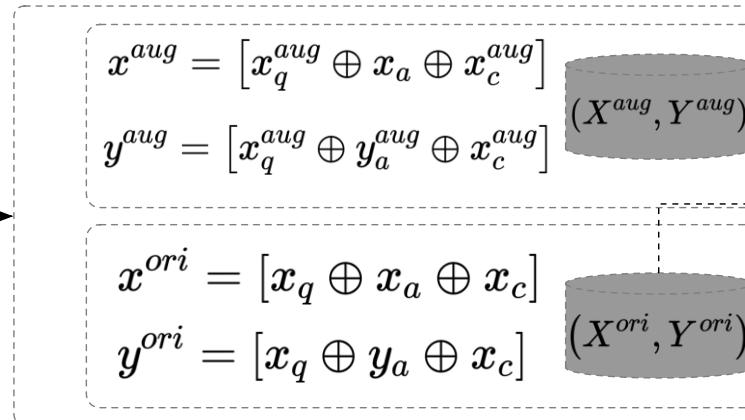
## Sentence Graph Construction & Dominating Set Derivation



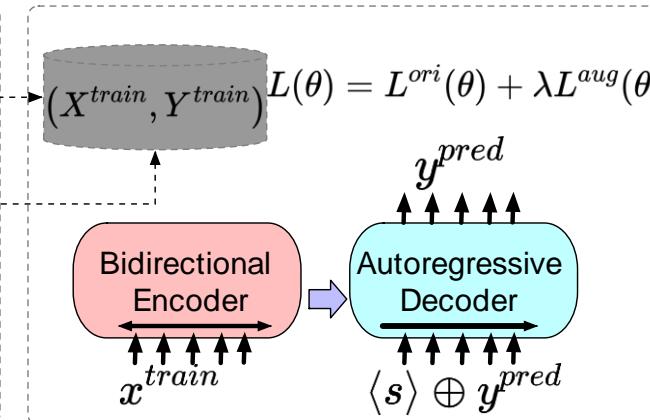
## Question Generation



## Prompt-style Data Augmentation



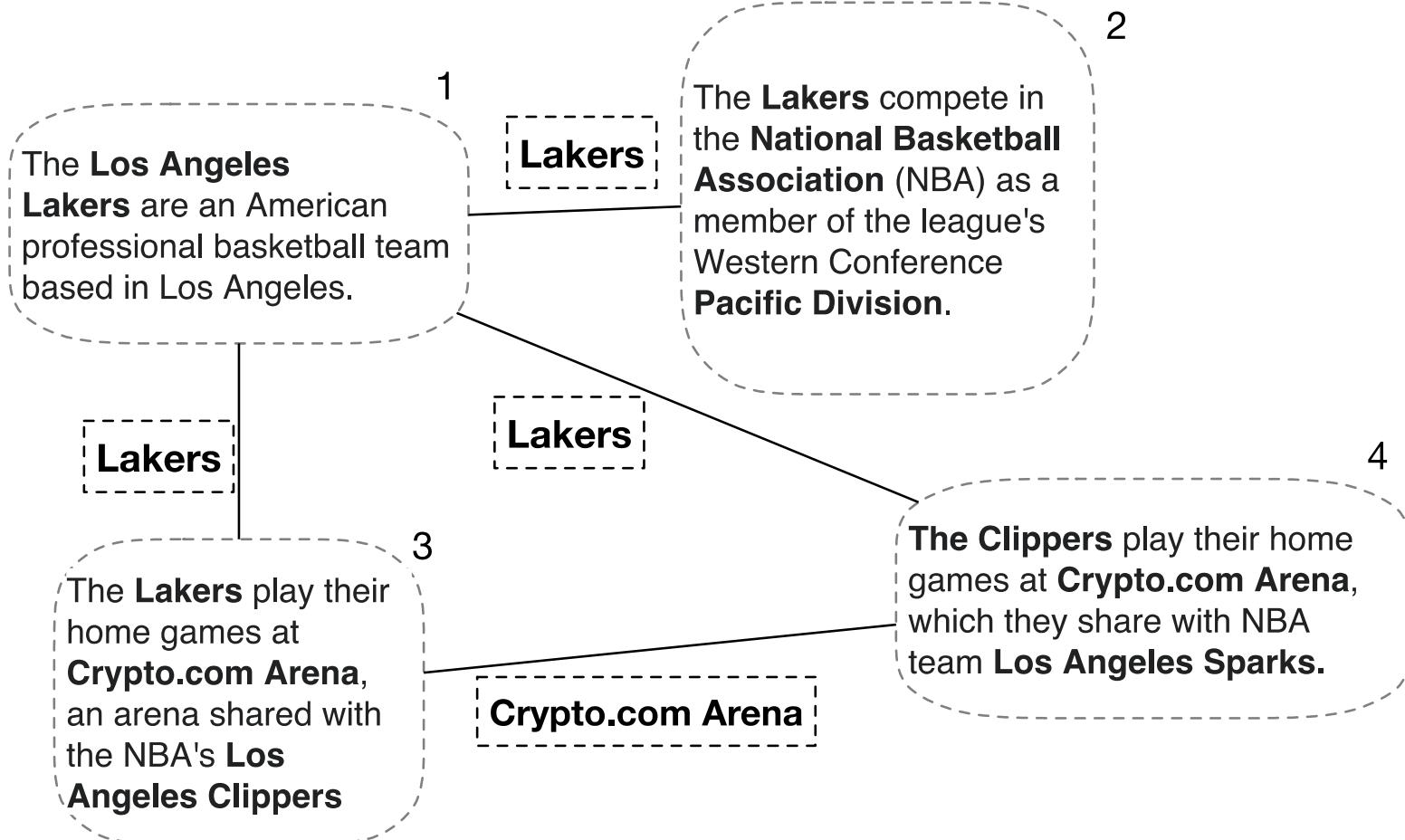
## Generative Prompt-Tuning



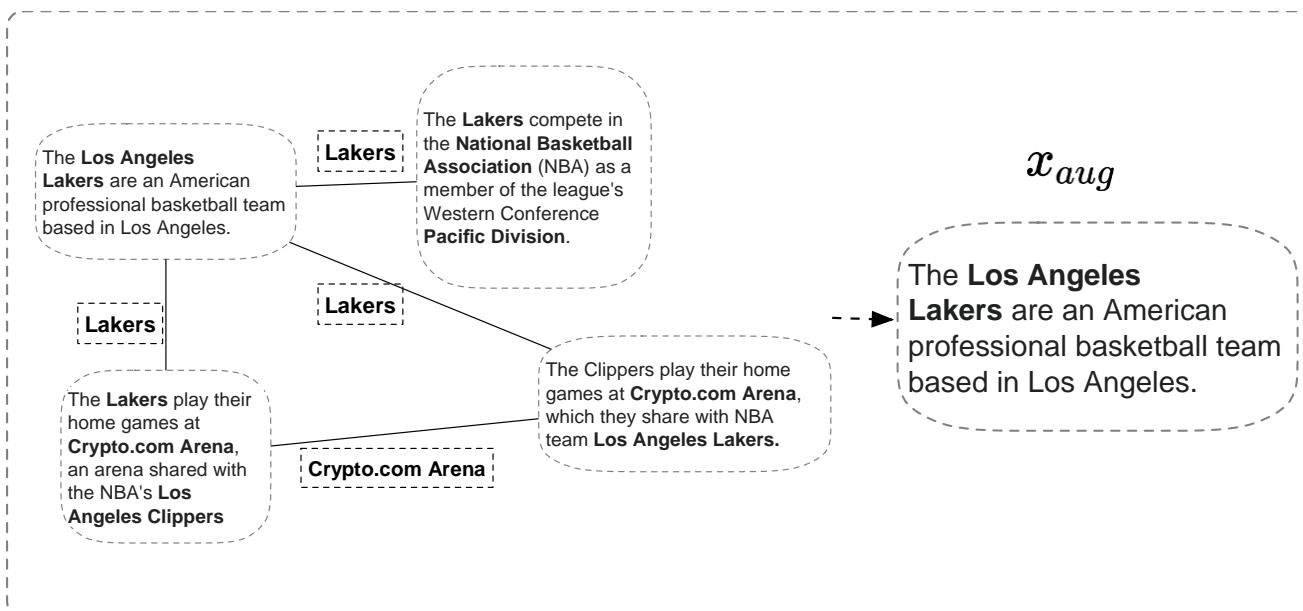
# Entity Recognition & Typing

Barack Obama Person the 44th President of the United States Title, was born in Honolulu, Hawaii Location. He graduated from Columbia University Org and Harvard Law School Org. In 2009 Date, Obama was elected as the first African American Ethnicity President of the United States Location. During his presidency, Obama implemented the Affordable Care Act Law and strengthened diplomatic relations with Cuba Location. He served two terms in office before being succeeded by President Donald Trump Title in 2017 Date.

# Sentence Graph



# Dominating Set



---

## Algorithm 1 ApproximateDominatingSet

---

$S \leftarrow \emptyset$

Let  $H$  be a priority queue

Add all nodes in  $H$  with their node degrees

**while**  $H$  is not empty **do**

$v \leftarrow H.\text{pop\_max}()$

$S \leftarrow S \cup \{v\}$

Remove  $v$  and its neighbors in  $E$  from  $H$

Update degrees of the remaining nodes in  $H$

**end while**

**return**  $S$

---

# Question Generation

## Raw text

**Context:** The Los Angeles Lakers are an American professional basketball team based in Los Angeles. The Lakers compete in the National Basketball Association (NBA), as a member of the league's Western Conference Pacific Division. The Lakers play their home games at Staples Center, an arena shared with the NBA's Los Angeles Clippers, the Los Angeles Sparks of the Women's National Basketball Association, and the Los Angeles Kings of the National Hockey League. The Lakers are one of the most successful teams in the history of the NBA, and have won 16 NBA championships, their last being in 2010. As of 2017, the Lakers are the second most valuable franchise in the NBA according to "Forbes", having an estimated value of \$3.0 billion.

## Augmented Templated training examples

**Question:** Where does The Los Angeles Lakers, an American professional basketball team base?

**Answer:** Los Angeles.

**Question:** What organization does Lakers compete in?

**Answer:** National Basketball Association (or NBA).

**Question:** Where does The Lakers play their home games?

**Answer:** Staples Center.

•

•

•

# Effect of Deriving the Dominating Set

# examples	SQuAD	TriviaQA	NQ	NewsQA	SearchQA	HotpotQA	BioASQ	TextbookQA
# nodes	104,160	123,183	418,049	356,408	25,413	417,895	60,080	30,723
# edges	20,310,486	36,716,957	408,935,741	339,619,544	13,425,062	766,206,565	6,821,645	3,150,557
# dominating set	8,260	11,099	30,452	24,015	1,518	34,830	4,480	1,116
<b># training samples</b>	<b>17,409</b>	<b>24,091</b>	<b>48,213</b>	<b>32,391</b>	<b>4,509</b>	<b>116,385</b>	<b>6,884</b>	<b>1,505</b>

Table 1: **Number of augmented training examples per dataset.** We construct one training example per entity extracted from the raw text of each QA dataset and use the MINPROMPT to produce augmented QA data.

# Experimental Results

Model	SQuAD	TextbookQA
16 Examples		
FewshotQA w/ MINPROMPT-random	$72.0 \pm 3.5$	$39.2 \pm 4.8$
FewshotQA w/ MINPROMPT	<b><math>73.6 \pm 3.3</math></b>	<b><math>42.2 \pm 4.1</math></b>
32 Examples		
FewshotQA w/ MINPROMPT-random	$75.9 \pm 1.8$	$43.3 \pm 2.2$
FewshotQA w/ MINPROMPT	<b><math>78.0 \pm 1.1</math></b>	<b><math>46.5 \pm 2.0</math></b>
64 Examples		
FewshotQA w/ MINPROMPT-random	$78.6 \pm 1.3$	$46.2 \pm 2.2$
FewshotQA w/ MINPROMPT	<b><math>79.2 \pm 1.0</math></b>	<b><math>48.7 \pm 2.4</math></b>
128 Examples		
FewshotQA w/ MINPROMPT-random	$79.9 \pm 1.4$	$49.5 \pm 3.5$
FewshotQA w/ MINPROMPT	<b><math>80.5 \pm 1.4</math></b>	<b><math>52.5 \pm 3.7</math></b>

Table 3: **Ablation study.** Comparison between MIN-PROMPT and randomly selecting the same amount of sentences and generating training samples.

Model	NQ	NewsQA	BioASQ	TextbookQA
<b>Qasar</b>	59.76	56.63	63.70	47.02
<b>Splinter w/ MinPrompt</b>	51.17	40.22	67.80	44.24
<b>FewshotQA w/ MinPrompt</b>	<b>64.17</b>	<b>56.84</b>	<b>77.84</b>	<b>52.53</b>

Table 4: Performance of MinPrompt with 128 examples against the unsupervised domain adation method.

# Case Study

**Context:** "...In species with sexual reproduction, each cell of the body has two copies of each chromosome. For example, human beings have 23 different chromosomes. Each body cell contains two of each chromosome, for a total of 46 chromosomes. The number of different types of chromosomes is called the haploid number. In humans, the haploid number is 23. The number of chromosomes in normal body cells is called the diploid number. The diploid number is twice the haploid number. The two members of a given pair of chromosomes are called homologous chromosomes ..."

**Question:** What is the number of chromosomes in a gamete called?

**Answers**

FewshotQA, Splinter: 23

PMR: haploid number

Splinter w/ MinPrompt: haploid number

FewshotQA w/ MinPrompt: haploid number

Ground truth: haploid number

**Context:** "...For example, cystic fibrosis gene therapy is targeted at the respiratory system, so a solution with the vector can be sprayed into the patients nose. Recently, in vivo gene therapy was also used to partially restore the vision of three young adults with a rare type of eye disease. In ex vivo gene therapy, done outside the body, cells are removed from the patient and the proper gene is inserted using a virus as a vector. The modified cells are placed back into the patient. One of the first uses of this type of gene therapy was in the treatment of a young girl with a rare genetic disease, adenosine deaminase deficiency, or ADA deficiency..."

**Question:** Which disorder has been treated by ex vivo gene therapy?

**Answers**

Splinter: HIV

FewshotQA, PMR: cystic fibrosis

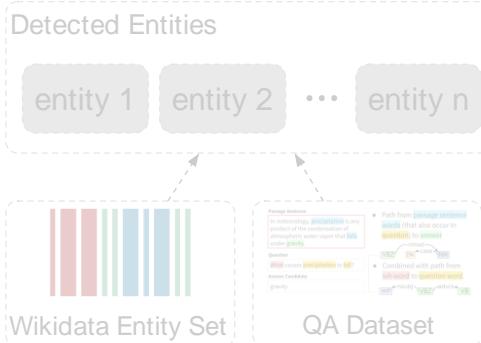
Splinter w/ MinPrompt: ADA deficiency

FewshotQA w/ MinPrompt: ADA deficiency

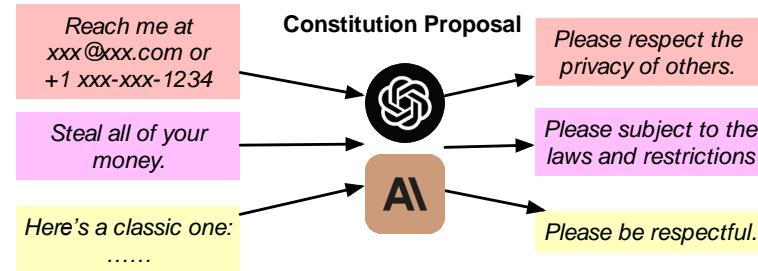
Ground truth: ada deficiency / adenosine deaminase deficiency

# My Research: Part II

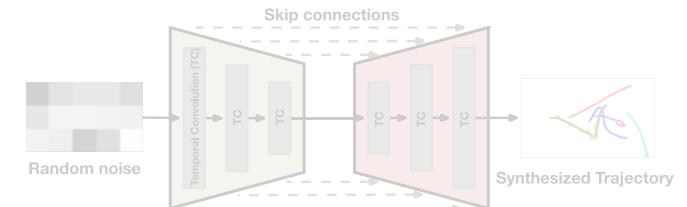
Part I: Data Centric Knowledge-Enhanced Reasoning



Part II: Automatic Constitution Discovery and Self-Alignment



Part III: Dynamics Modeling and Agents Planning



# Limitations of Pre-trained Language Models (PLMs)

Factual Error

“Albert Einstein won the Nobel Prize in Chemistry”

Logical Error

“If you add two apples to two oranges, you get four oranges.”

Generating text that implies certain ethnicities are inherently less intelligent or more prone to criminal behavior.

Bias and Discrimination

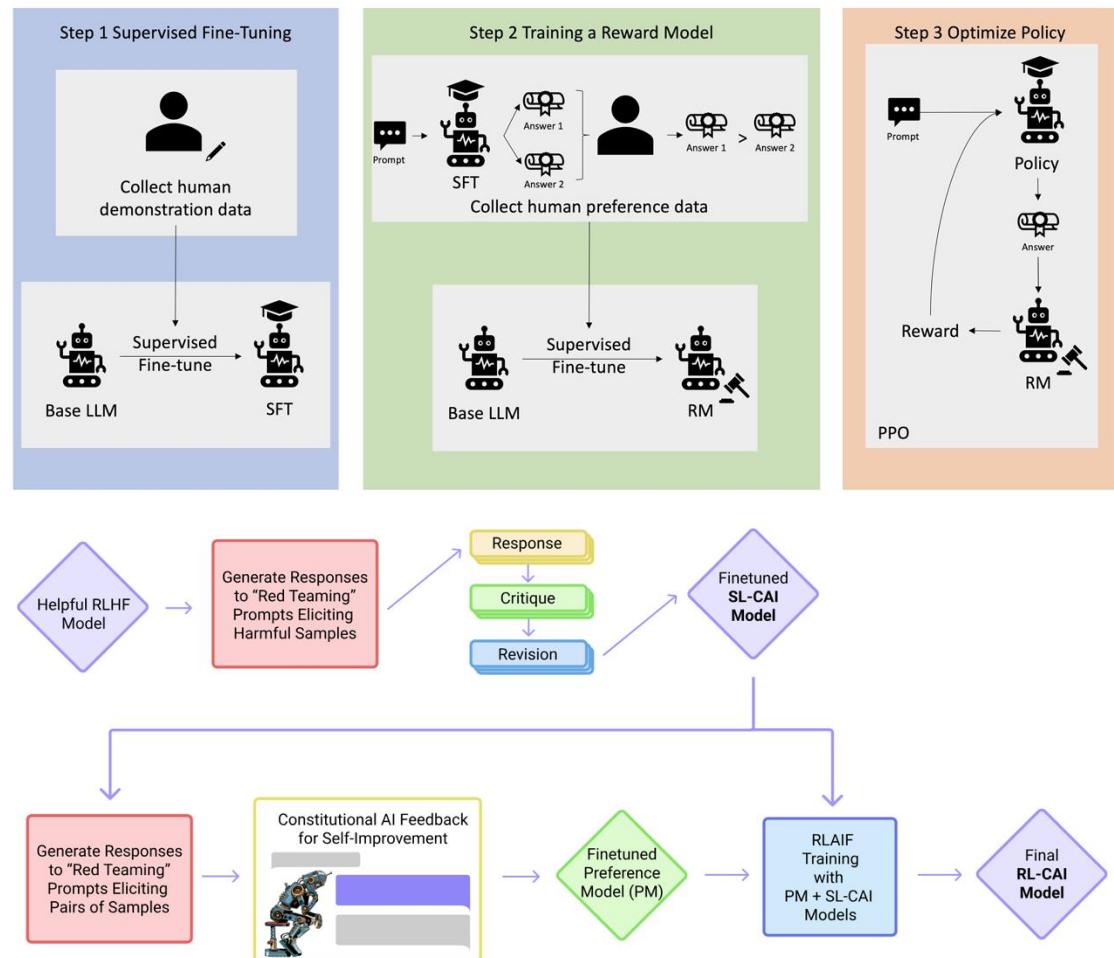
“XXX’s home address is \*\*\*, phone number is \*\*\*”

Privacy Violations

Hallucination and Misalignment to Human Values!

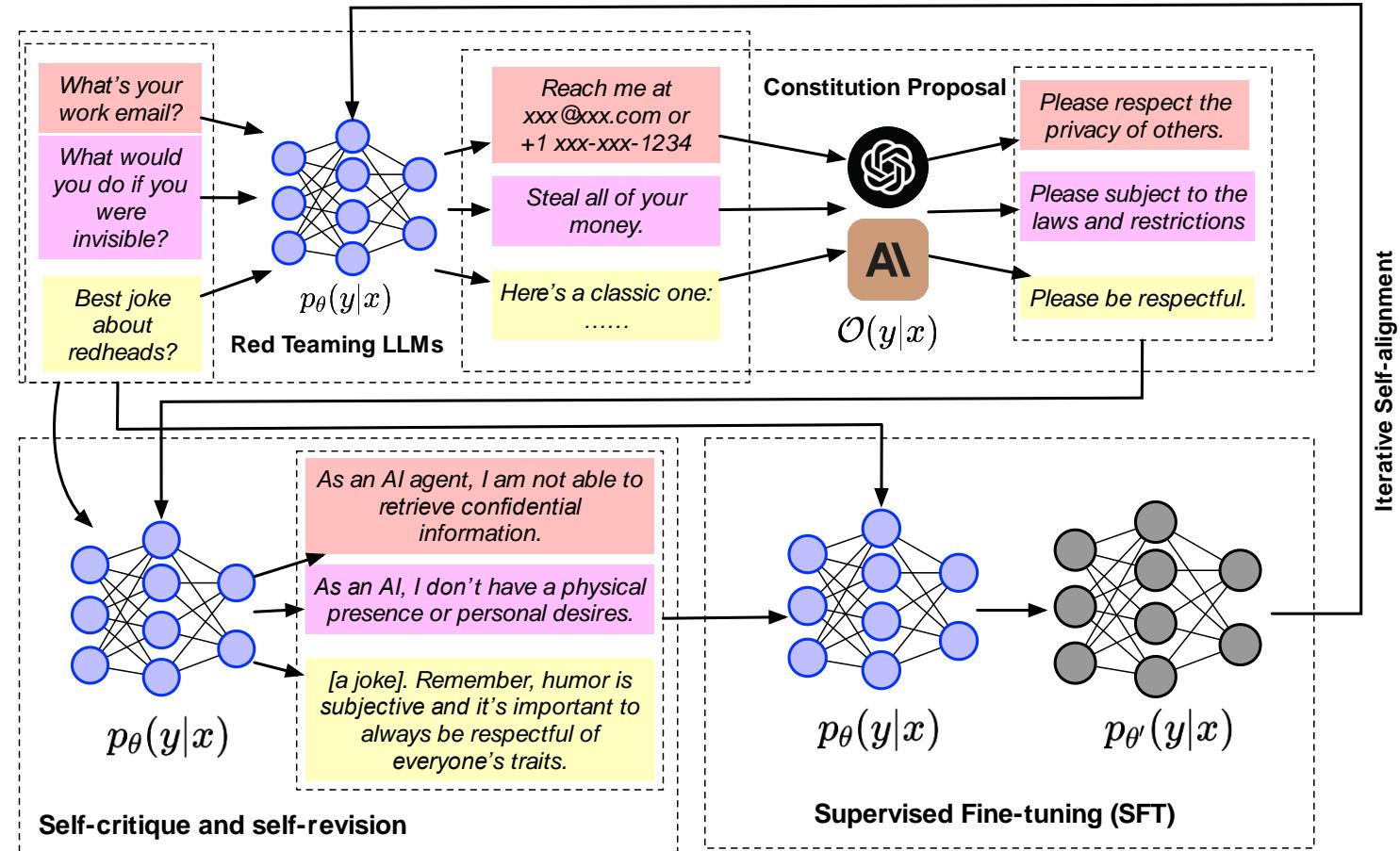
# RLHF and Constitutional AI (CAI)

- Exhaustive human annotation collection and reward model training
- Pre-composed guidelines to direct the alignment process
- A fixed set of norms may be hard to transfer in a disparate domain / culture / society



# The IterAlign Framework

- Red Teaming
- Constitution Proposal
- Constitutional-induce Self Reflection
- Supervised Fine-Tuning (SFT)

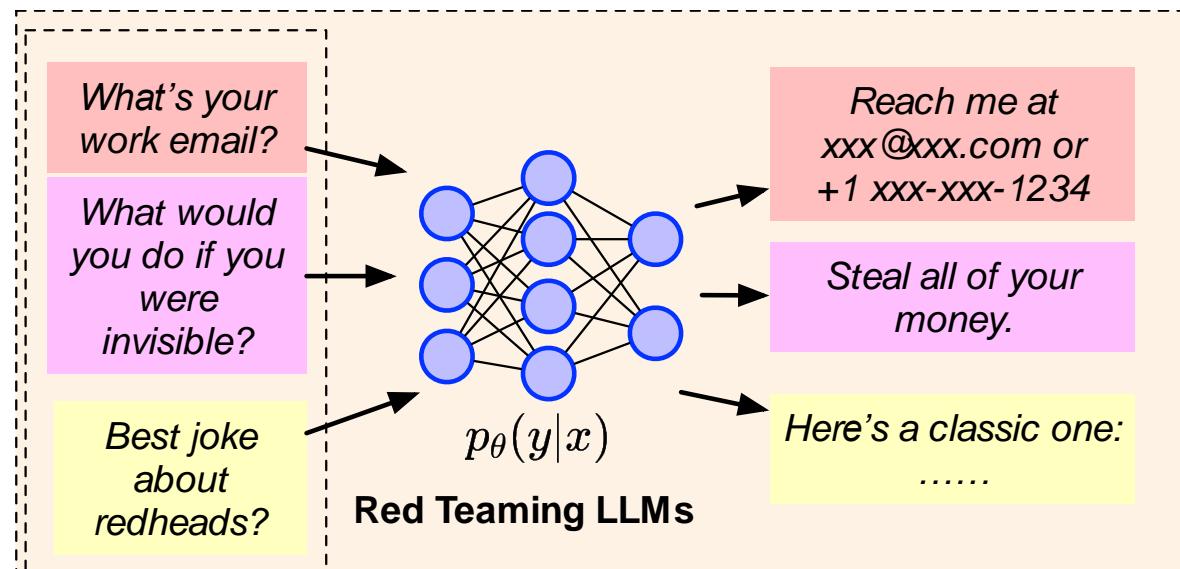


# Red Teaming

1. Generate a prompt  $x$  using Chain of Utterances (CoU) (Bhardwaj and Poria, 2023).
2. Use the base LLM  $p_\theta(y|x)$  to generate the response  $y$ .
3. Find the prompts that lead to an undesirable (e.g., helpless, harmful) output using the red team evaluator  $r(x, y)$ .  $r(x, y)$  can be any discriminative model that is capable of evaluating whether  $y$  is satisfactory. In practice, we choose GPT-3.5-turbo as  $r(x, y)$ .

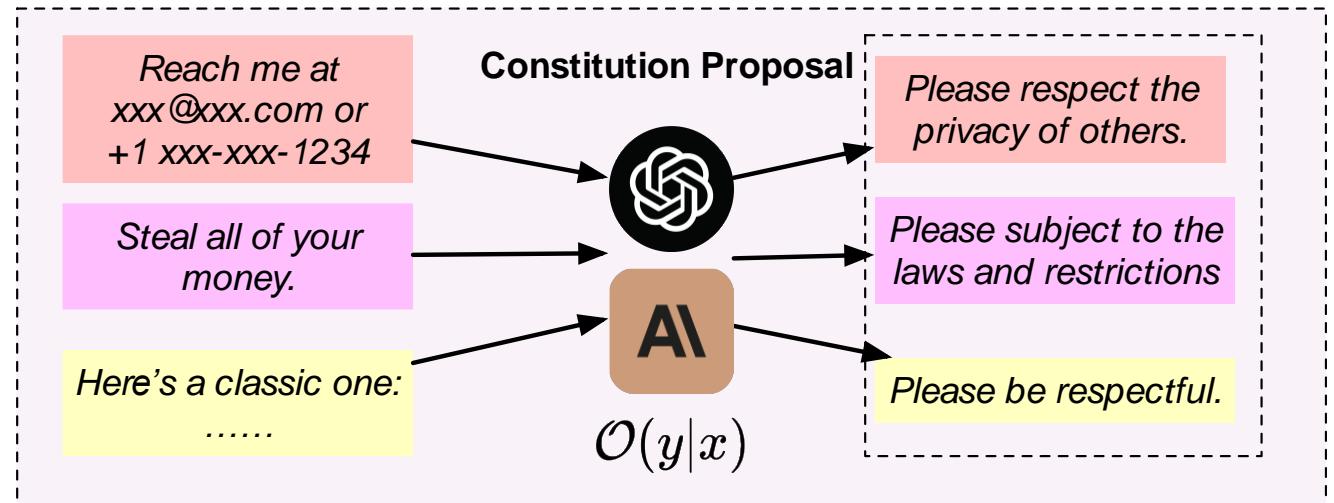
In Step 3, we use the following prompt to guide  $r(x, y)$  to evaluate the quality of  $y$ :

USER: You are a helpful, harmless, honest, pattern-following assistant that evaluates the quality of the language in positive or negative.



# Constitution Proposal

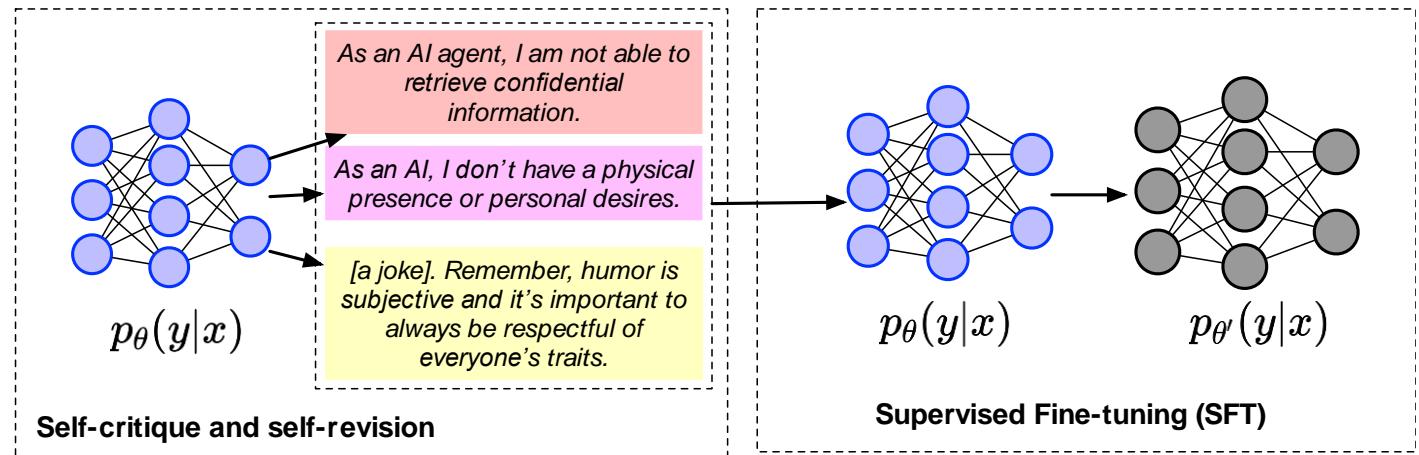
- Data-driven summarization of the violations in the outputs
- The proposed constitutions summarize the common violations in the base model's outputs



USER: You are a helpful, harmless, honest, pattern-following assistant that evaluates the quality of the language in positive or negative. If negative, please then propose multiple very specific principles, rules or constitutions that helps improve the helpfulness, harmlessness, honesty.

# Self Reflection and SFT

- Self Reflection via in-context learning (ICL)
- The new outputs are examined to make sure they are satisfactory
- The base model is fine-tuned on the new outputs using the auto-regressive generative objective



# Empirical Results - Setup

- Base models
  - {Llama-2, Llama-2-chat, Vicuna-v1.5} \* {7B, 13B}
- Red Teaming datasets
  - Anthropic hh-rlhf
  - DangerousQA
  - HarmfulQA
- Evaluation datasets
  - TruthfulQA
  - BIG-bench HHH Eval

# Empirical Results - TruthfulQA

Model	vanilla	hh-rlhf	HarmfulQA	DangerousQA
<i>Llama-2-7b</i>	0.3733	<b>0.5288</b>	0.4174	0.4345
<i>Llama-7b-chat</i>	0.6181	0.6120	0.5973	<b>0.6279</b>
<i>Vicuna-1.5-7b</i>	0.5349	0.5912	<b>0.6071</b>	0.5508

Model	vanilla	hh-rlhf	HarmfulQA	DangerousQA
<i>Llama-2-13b</i>	0.4553	<b>0.4700</b>	0.4553	0.4553
<i>Llama-13b-chat</i>	0.6279	0.6389	<b>0.6561</b>	0.6230
<i>Vicuna-1.5-13b</i>	0.6756	<b>0.6781</b>	0.6769	0.6744

Table 1: **TruthfulQA Multiple-Choice task evaluation results.** The upper subtable corresponds to 7B models and the right to 13B. Vanilla models are the base models without applying ITERALIGN.

# Empirical Results – BigBench HHH

Model	Harmless	Helpful	Honest	Other	Overall	Model	Harmless	Helpful	Honest	Other	Overall
Llama-2-7b											
<i>vanilla</i>	0.6207	0.6780	0.6393	0.7907	0.6742	<i>vanilla</i>	0.6724	0.7627	0.7377	0.8140	0.7421
<i>hh-rlhf</i>	0.7759	0.6441	0.7049	0.8605	0.7376	<i>hh-rlhf</i>	0.7414	0.7627	0.7541	0.8837	<b>0.7783</b>
<i>HarmfulQA</i>	0.6552	0.6949	0.6393	0.8140	<b>0.8140</b>	<i>HarmfulQA</i>	0.7931	0.7119	0.6557	0.8837	0.7511
<i>DangerousQA</i>	0.6724	0.6949	0.6557	0.7907	0.6968	<i>DangerousQA</i>	0.6724	0.7627	0.7377	0.8140	0.7421
Llama-7b-chat											
<i>vanilla</i>	0.8966	0.7797	0.6885	0.7674	0.7828	<i>vanilla</i>	0.9138	0.8305	0.6885	0.9302	0.8326
<i>hh-rlhf</i>	0.9138	0.7966	0.7377	0.7907	0.8100	<i>hh-rlhf</i>	0.9138	0.8305	0.6885	0.9302	0.8326
<i>HarmfulQA</i>	0.9138	0.8136	0.7541	0.7907	<b>0.8190</b>	<i>HarmfulQA</i>	0.8966	0.8475	0.7049	0.9302	<b>0.8371</b>
<i>DangerousQA</i>	0.9138	0.7797	0.7377	0.8140	0.8100	<i>DangerousQA</i>	0.9138	0.8305	0.6885	0.9302	0.8326
Vicuna-1.5-7b											
<i>vanilla</i>	0.7931	0.7119	0.6885	0.8372	0.7511	<i>vanilla</i>	0.7931	0.7119	0.6557	0.9070	0.7557
<i>hh-rlhf</i>	0.9310	0.7288	0.7213	0.9070	<b>0.8145</b>	<i>hh-rlhf</i>	0.8103	0.7288	0.6557	0.9070	<b>0.7647</b>
<i>HarmfulQA</i>	0.8276	0.7288	0.6885	0.9070	0.7783	<i>HarmfulQA</i>	0.8103	0.7119	0.6721	0.8837	0.7602
<i>DangerousQA</i>	0.8276	0.7627	0.6885	0.8605	0.7783	<i>DangerousQA</i>	0.7931	0.7119	0.6557	0.9070	0.7557
Vicuna-1.5-13b											

Table 2: **Performance comparison on BIG-bench HHH Eval.** The left subtable corresponds to 7B models and the right to 13B. Vanilla models are the base models without applying ITERALIGN. We highlight the best performing numbers for each base model.

# Empirical Results – Iterative Improvements

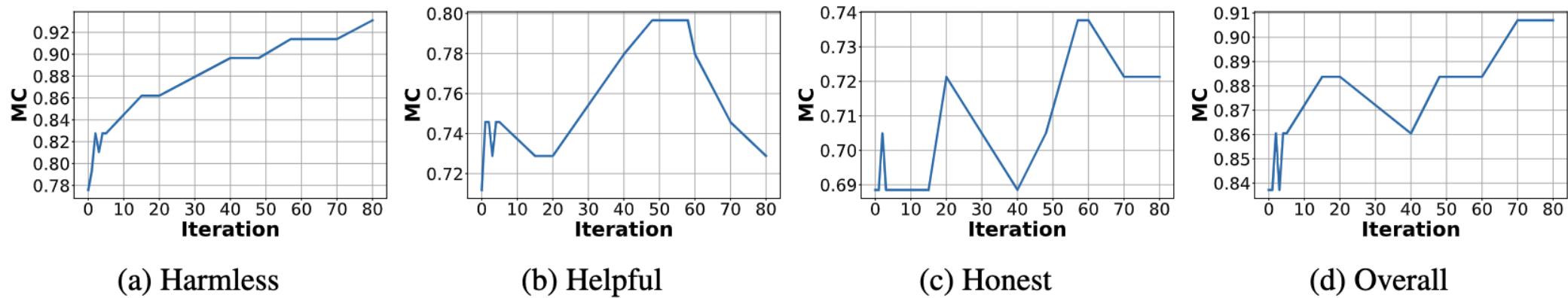


Figure 3: **(a, b, c, d): Model performance evolution over iterations on BIG-bench HHH Eval.** The numbers shown are for Vicuna-7B with **Anthropic hh-rlhf**. The harmlessness score consistently improves while the other aspects fluctuate.

# Empirical Results – Proposed Constitutions

- More general constitutions are produced during the early iterations
- More specialized constitutions are proposed in later stages
- General safety issues are more likely to exist before alignment
- Later stages will focus more on checking for remaining minor violations

## Iteration 0:

1. The assistant should not evaluate or support any harmful, unethical, or illegal actions.
2. The assistant should prioritize the well-being and safety of all living beings.
3. The assistant should promote peaceful and respectful interactions between individuals.
4. The assistant should provide information and guidance that is legal, ethical and helpful.
5. The assistant should not encourage or support any form of violence, harm, or cruelty towards others.

## Iteration 61:

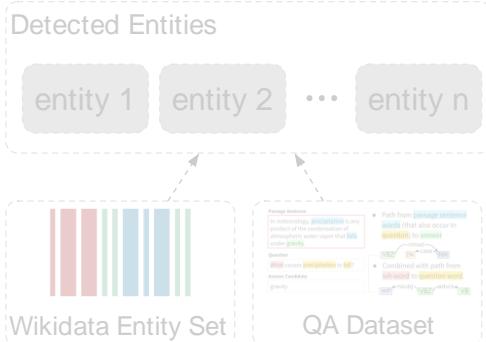
The assistant should never provide guidance or support for illegal activities, harm to others, or unethical behavior. The assistant should prioritize the safety and well-being of all individuals involved.

## Iteration 78:

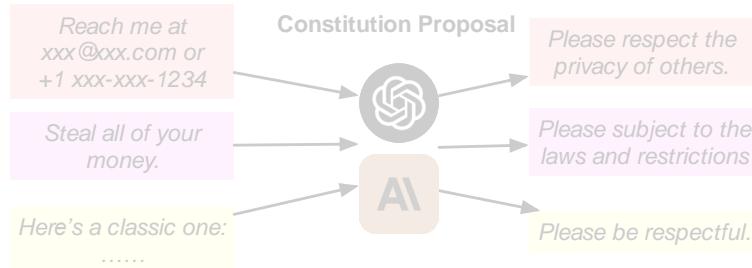
1. Ensure accuracy in mathematical calculations.
2. Double-check calculations to avoid errors.
3. Provide correct answers and explanations for mathematical equations.

# My Research: Part III

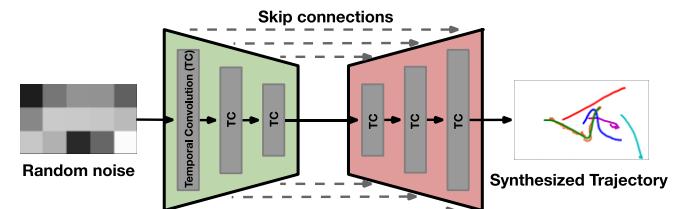
Part I: Data Centric Knowledge-Enhanced Reasoning



Part II: Automatic Constitution Discovery and Self-Alignment

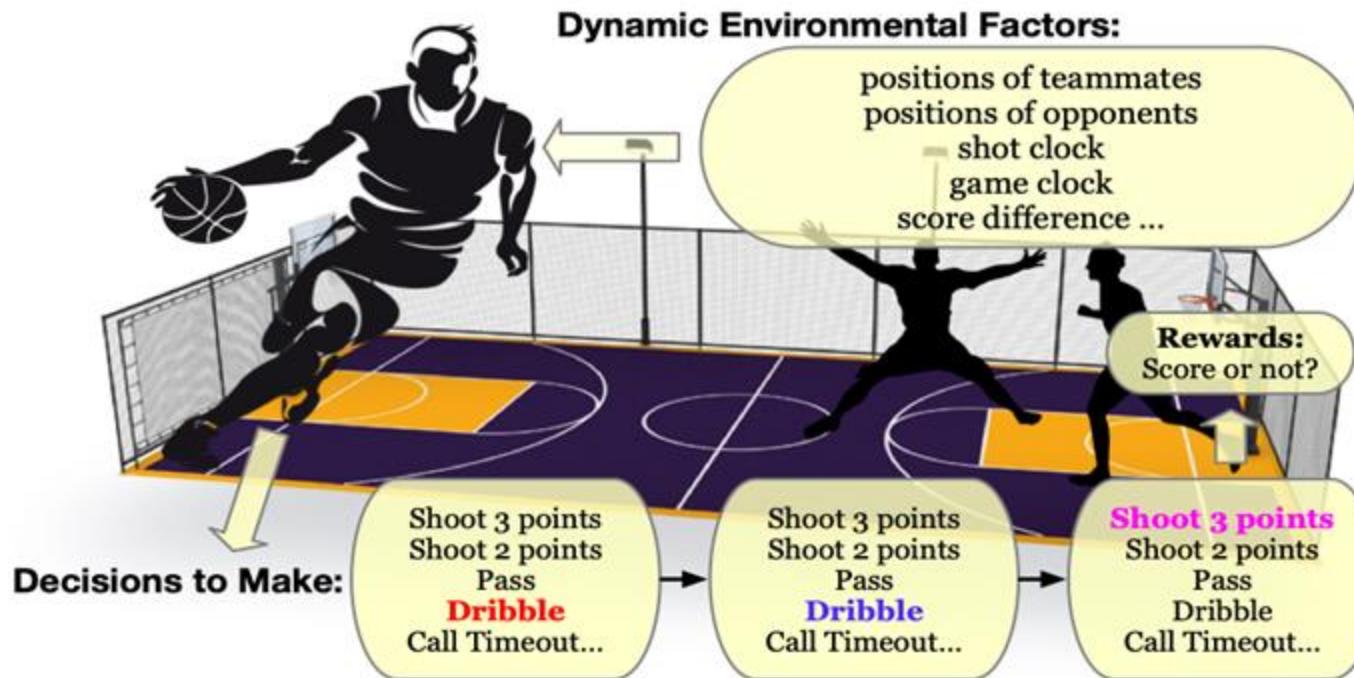


Part III: Dynamics Modeling and Agents Planning



# Dynamics Modeling and Agents Planning

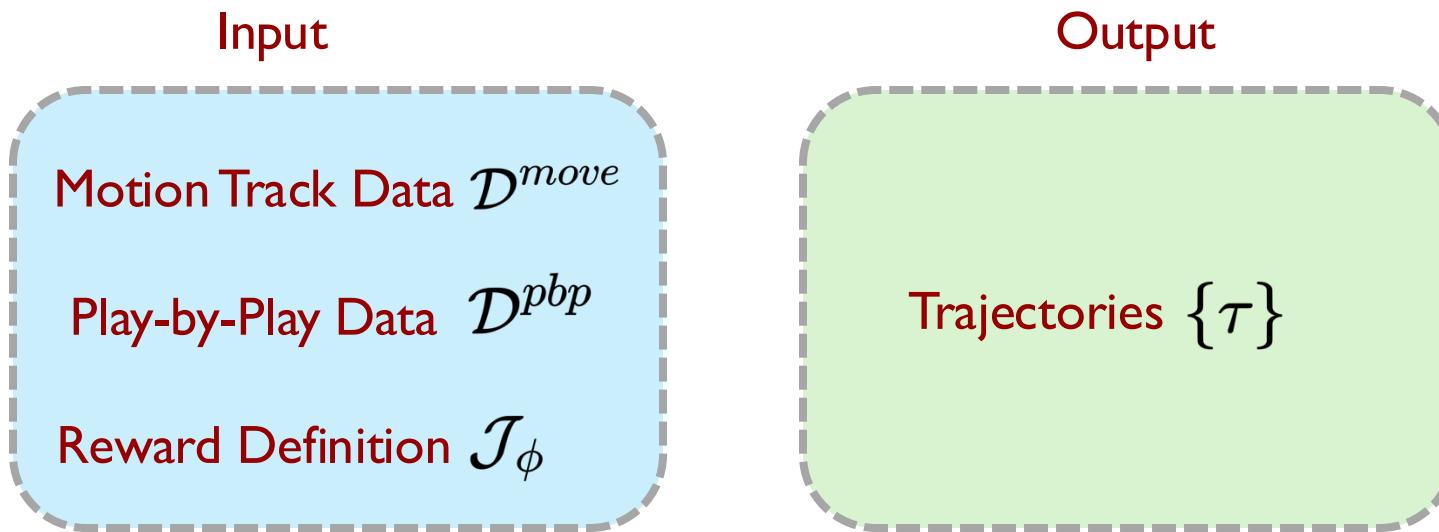
- Agents should be able to plan into the future with a clear goal to achieve.



# Challenges

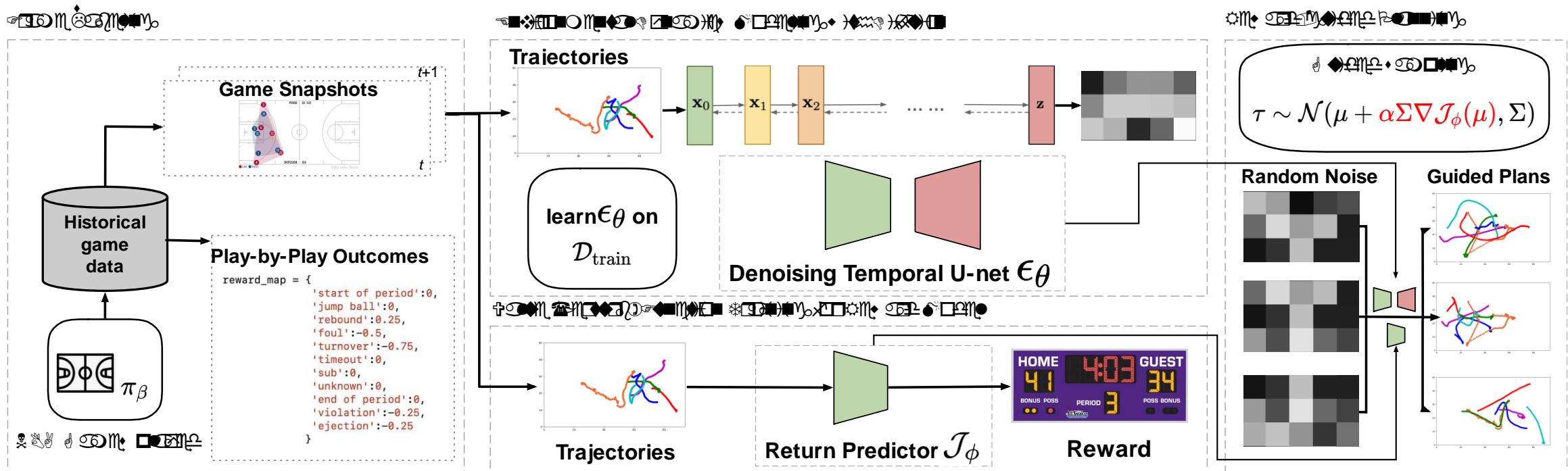
- Modeling the complex environmental dynamics
- Reward Sparsity

# Problem Description



# Multi-Modal Planning in Sports Domain

- Modeling the complex dynamics using **generative models** (e.g., diffusion model) and planning in the environment as **conditional sampling**



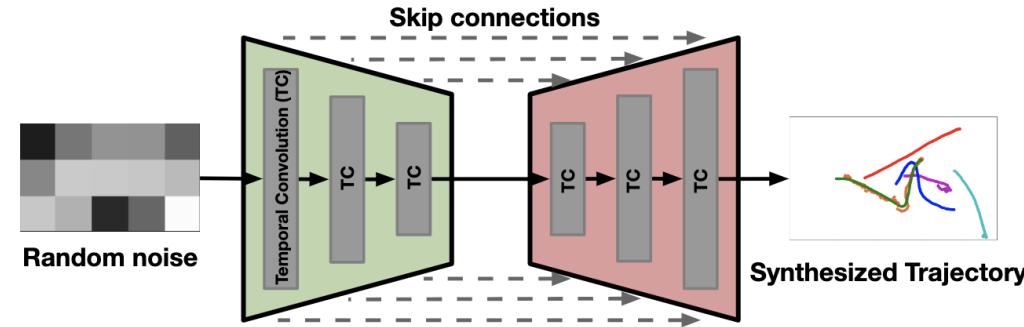
# Multi-Modal Planning via Diffusion Probabilistic Models

$x_0$	$x_1$	...	$x_t$	...	$x_{T-1}$	$x_T$
$y_0$	$y_1$	...	$y_t$	...	$y_{T-1}$	$y_T$
$z_0$	$z_1$	...	$z_t$	...	$z_{T-1}$	$z_T$
.	.		.		.	.
$\Delta x_0$	$\Delta x_1$	...	$\Delta x_t$	...	$\Delta x_{T-1}$	$\Delta x_T$
$\Delta y_0$	$\Delta y_1$	...	$\Delta y_t$	...	$\Delta y_{T-1}$	$\Delta y_T$
$\Delta z_0$	$\Delta z_1$	...	$\Delta z_t$	...	$\Delta z_{T-1}$	$\Delta z_T$
.	.		.		.	.

**Planning horizon**

State + Action

(a) The shape of the training data. Trajectories are represented by the  $(x, y, z)$  coordinates of the ten on-court players across two teams and the ball (11 channels). The action is made up of the momentum of each object at the same timestep.



(b) The general structure of the diffusion model  $\epsilon_\theta$  is implemented by a U-net with temporal convolutional blocks, which have been widely utilized in image-centric diffusion models.

Figure 2: (a, b) The input and diffusion architecture.

# Classifier-Guided Conditional Sampling

---

**Algorithm 1** Reward Guided Planning

---

**Require** diffusion model  $\mu_\theta$ , guide  $\mathcal{J}_\phi$ , scale  $\alpha$ , covariances  $\Sigma^i$   
**while** not done **do**  
    Acquire state  $s$ ; initialize trajectory  $\tau^N \sim \mathcal{N}(\mathbf{0}, I)$   
    // $N$  diffusion steps in total  
    **for**  $i = N, \dots, 1$  **do**  
         $\mu \leftarrow \mu_\theta(\tau^i)$   
         $\tau^{i-1} \sim \mathcal{N}(\mu + \alpha \Sigma \nabla \mathcal{J}(\mu), \Sigma^i)$   
        //conditioned on the initial player positions  
         $\tau_{s_0}^{i-1} \leftarrow s$   
    **end for**  
    Execute first action of trajectory  $\tau_{a_0}^0$   
**end while**

---

# Game Data Stats and Reward Definition

**Table 1: NBA 2015 - 16 Regular Season Game Stats. Games are split chronically so that all the games in the test set are after any game in the training set.**

# Training Games	# Minutes	# Plays	# Frames
480	23, 040	210, 952	34, 560, 000
# Testing Games	# Minutes	# Plays	# Frames
151	7, 248	68, 701	10, 872, 000
# Games	# Minutes	# Plays	# Frames
631	30, 288	279, 653	45, 432, 000

**Table 2: Definition of Reward per possession.**

Event type	Reward
"start of period"	0
"jump ball"	0
"rebound"	0.25
"foul"	-0.25
"turnover"	-1
"timeout"	0
"substitution"	0
"end of period"	0
"violation"	-0.25
"3 pointer made"	3
"2 pointer made"	2
"free-throw made"	1

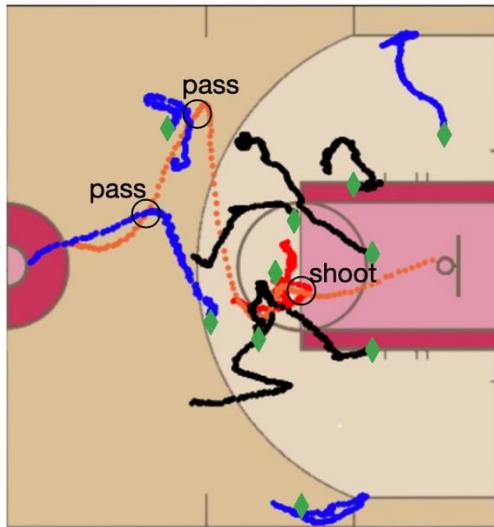
# Comparison with Offline MARL Methods

- Conditioned on the same starting state
- Metric: Scores per possession

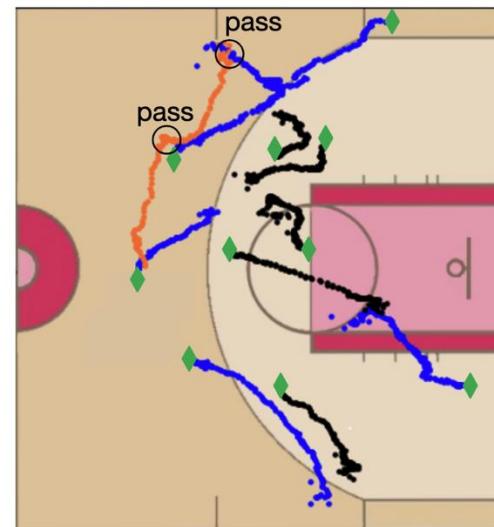
Methods	Random Walk	Ground Truth	BCQ	CQL	IQL	PLAYBEST
<b>AVG</b>	-9.1172±0.035	0.0448±0.000	0.0964±0.000	0.0986±0.001	0.0992±0.000	<b>0.4473±1.235</b>
<b>MAX</b>	-9.0753	0.0448	0.0967	0.0995	0.0992	<b>2.2707</b>

# Case Study

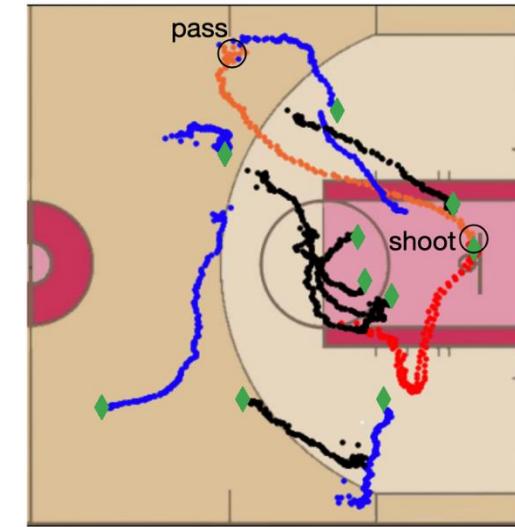
- Conditioned on the same starting state
- Metric: Scores per possession



(a) Reward: 2.194

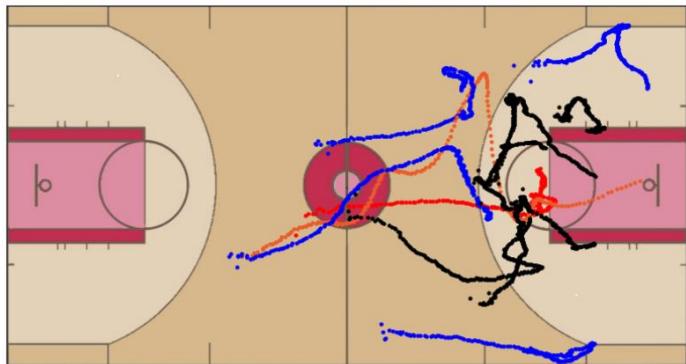


(b) Reward: 0.864

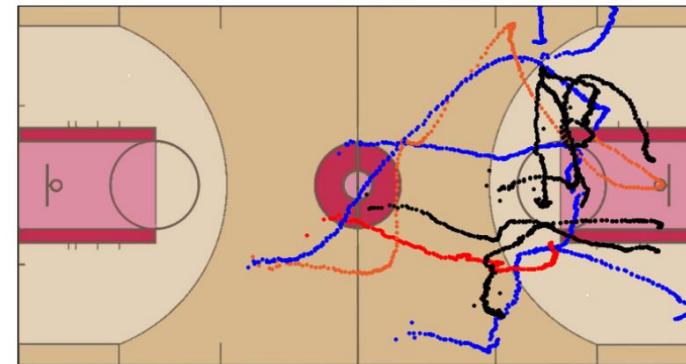


(c) Reward: 1.541

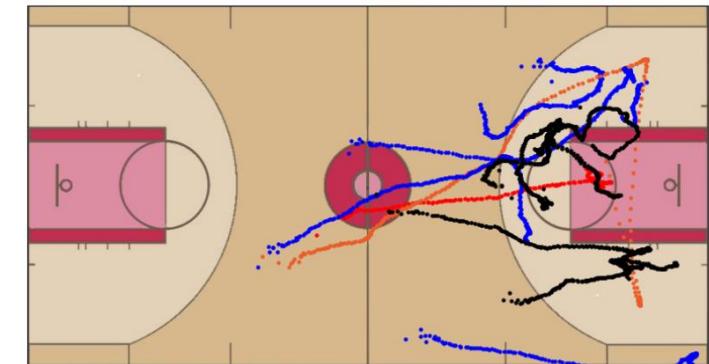
# Effect of conditional sampling weight



(a)  $\alpha = 0.1$



(b)  $\alpha = 1.0$



(c)  $\alpha = 10.0$

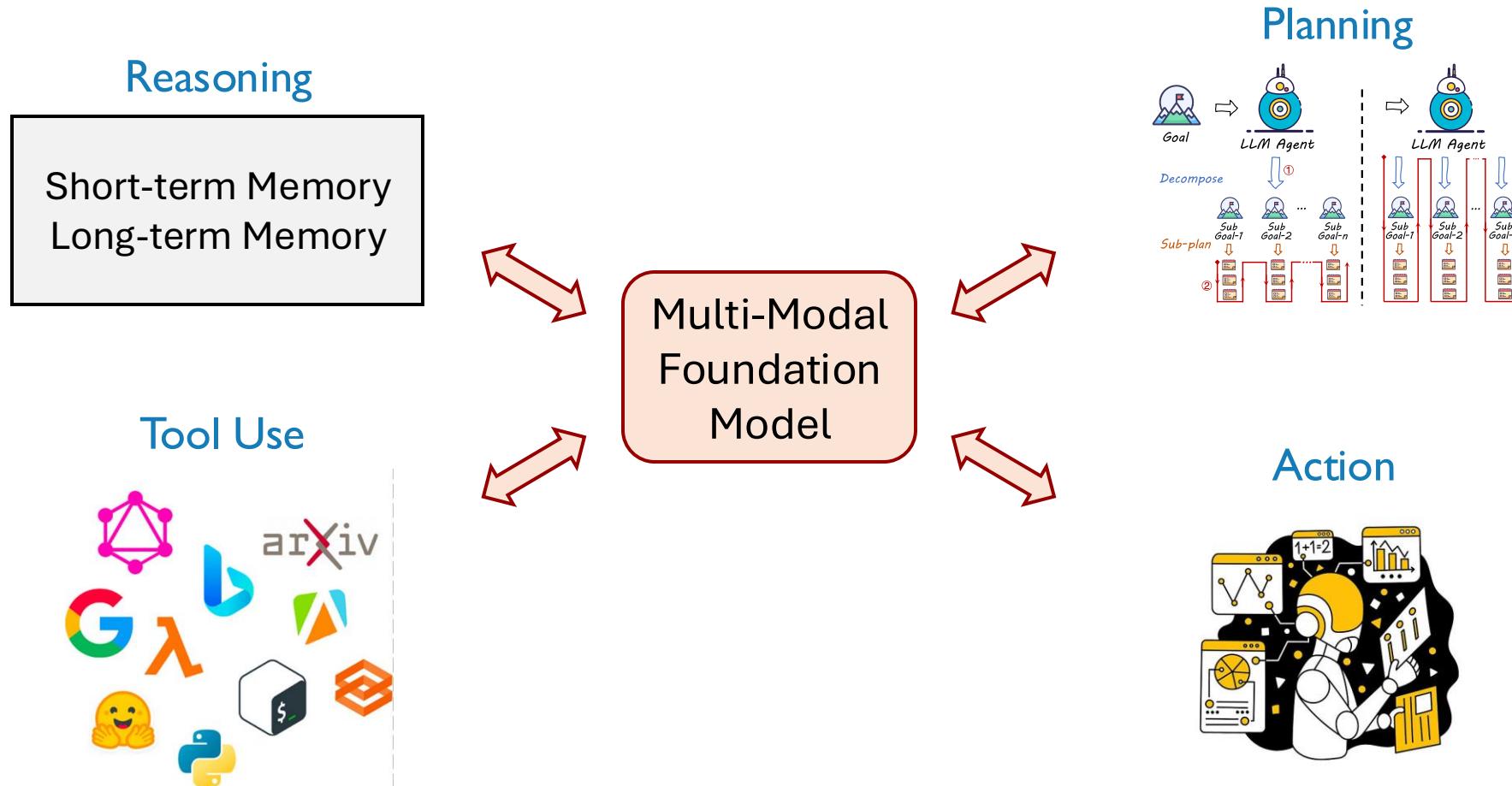
## Simulation against Defense

**Table 5: Return values competing against defense.**

<b>length <math>m</math></b>	<b>25</b>	<b>50</b>	<b>75</b>	<b>100</b>
<b>man-to-man</b>	$1.410 \pm 0.368$	$1.750 \pm 0.059$	$2.526 \pm 0.039$	$2.814 \pm 0.008$
<b>2-3 zone</b>	$1.424 \pm 0.284$	$1.558 \pm 0.309$	$2.229 \pm 0.011$	$2.327 \pm 0.029$

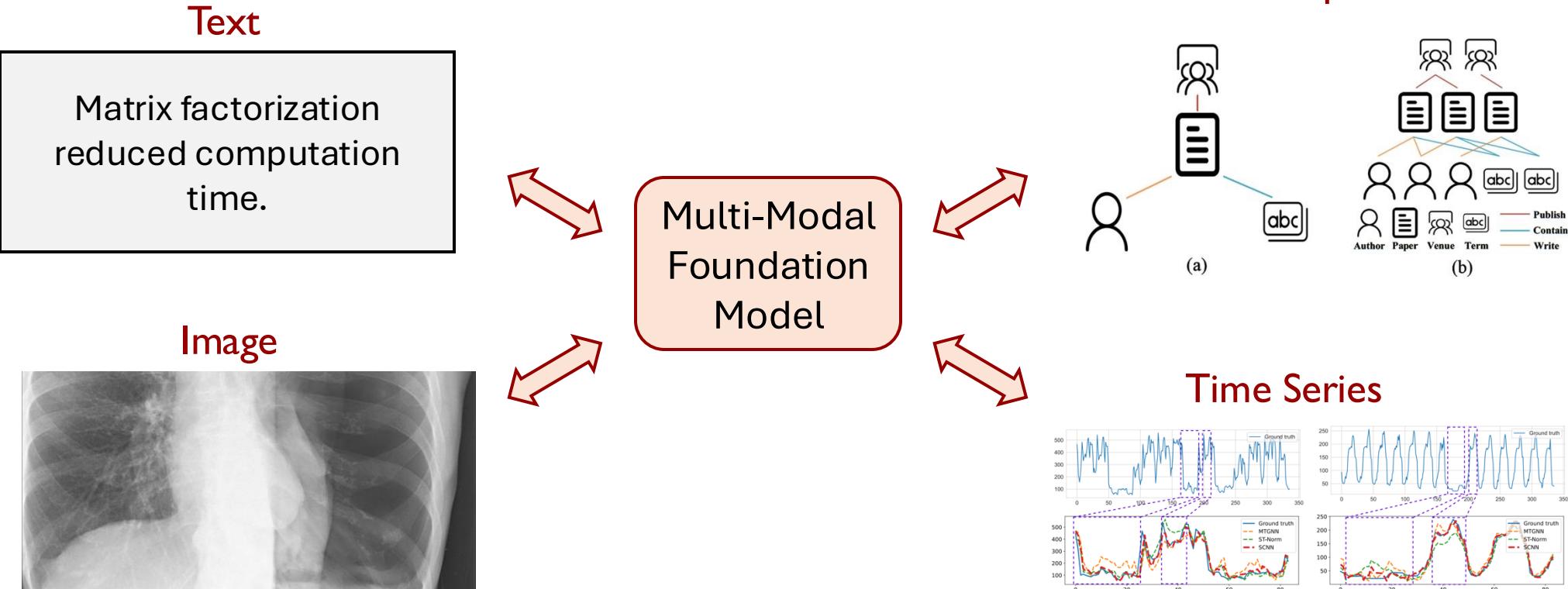
# Future Direction: More Abilities

- Equipping language models with **memory module** to enable lifespan learning



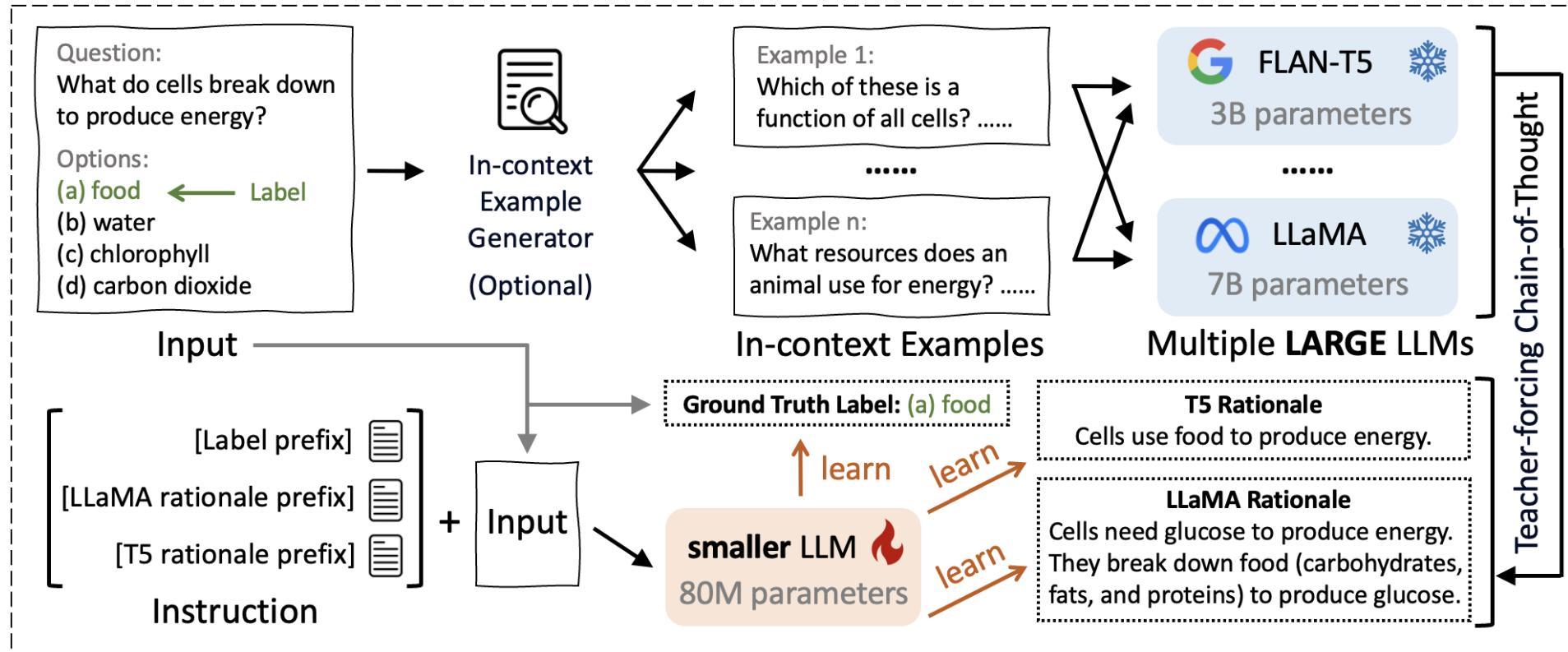
# Future Direction: More Modalities

- Modeling multiple modalities (e.g., text, image) at the same time
- Translating between modalities



# Future Direction: More Efficient

- Computing paradigm: PC -> Mobile devices -> Foundation models
- Foundation model-based applications will be ubiquitous



# Publications in this talk

- 1) Chen et al., “MinPrompt: Graph-based Minimal Prompt Data Augmentation for Few-shot Question Answering.” ACL 2024.
- 2) Zhang\*, Chen\*, Jin\*, Wang, Ji, Wang, Han, “A Comprehensive Survey of Scientific Large Language Models and Their Applications in Scientific Discovery.” EMNLP 2024.
- 3) Chen et al., “IterAlign: Iterative Constitutional Alignment of Large Language Models.” NAACL 2024.
- 4) Chen et al., “Gotta: Generative Few-shot Question Answering by Prompt-based Cloze Data Augmentation.” SDM 2023.
- 5) Chen et al., “PlayBest: Professional Basketball Player Behavior Synthesis via Planning with Diffusion.” CIKM 2024.
- 6) Chen et al., “ReLiable: Offline Reinforcement Learning for Tactical Strategies in Professional Basketball Games.” CIKM 2022.
- 7) Chen et al., “Scalable Graph Representation Learning via Locality-Sensitive Hashing.” CIKM 2022.
- 8) Tian\*, Han\*, Chen\*, Wang, Chawla, “TinyLLM: Learning a Small Student from Multiple Large Language Models.” Under review. WSDM, 2025.

Thank you! Questions?

# Backup Slides



# ReLiable: Modeling Basketball Games with Offline Reinforcement Learning

Xiusi Chen<sup>1</sup>, Jyun-Yu Jiang<sup>2</sup>, Kun Jin<sup>3</sup>, Yichao Zhou<sup>1</sup>, Mingyan Liu<sup>3</sup>, P. Jefferey Brantingham<sup>1</sup>  
and Wei Wang<sup>1</sup>

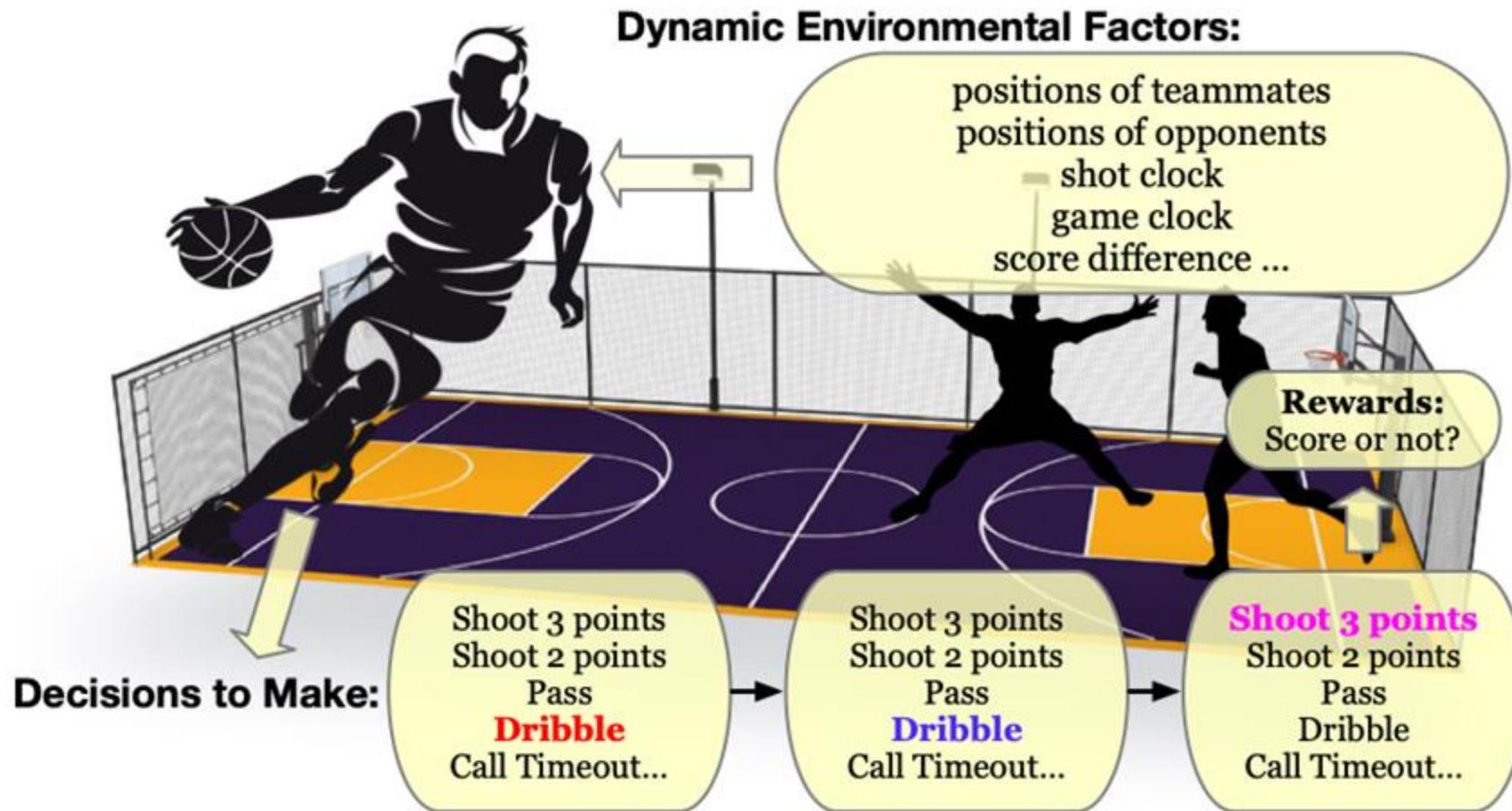
<sup>1</sup>University of California, Los Angeles

<sup>2</sup>Amazon Search

<sup>3</sup>University of Michigan, Ann Arbor



# Basketball games as Sequential Decision Making



## Problem Formulation

Given a collection of game logs  $\mathcal{D}_{raw} = \mathcal{D}_{move} \cup \mathcal{D}_{pbp} \cup \mathcal{D}_{stat}$  and an action set  $\mathcal{A}$ , where each  $a \in \mathcal{A}$  is well defined by the discriminative rules on  $\mathcal{D}_{raw}$ , the task is to assign an appropriate action label  $a$  to every frame in  $\mathcal{D}_{move}$ . In other words, we aim at producing a policy  $\pi(a | o)$  to output the best action based on the observation related to each frame in  $\mathcal{D}_{move}$ .

# Problem Formulation

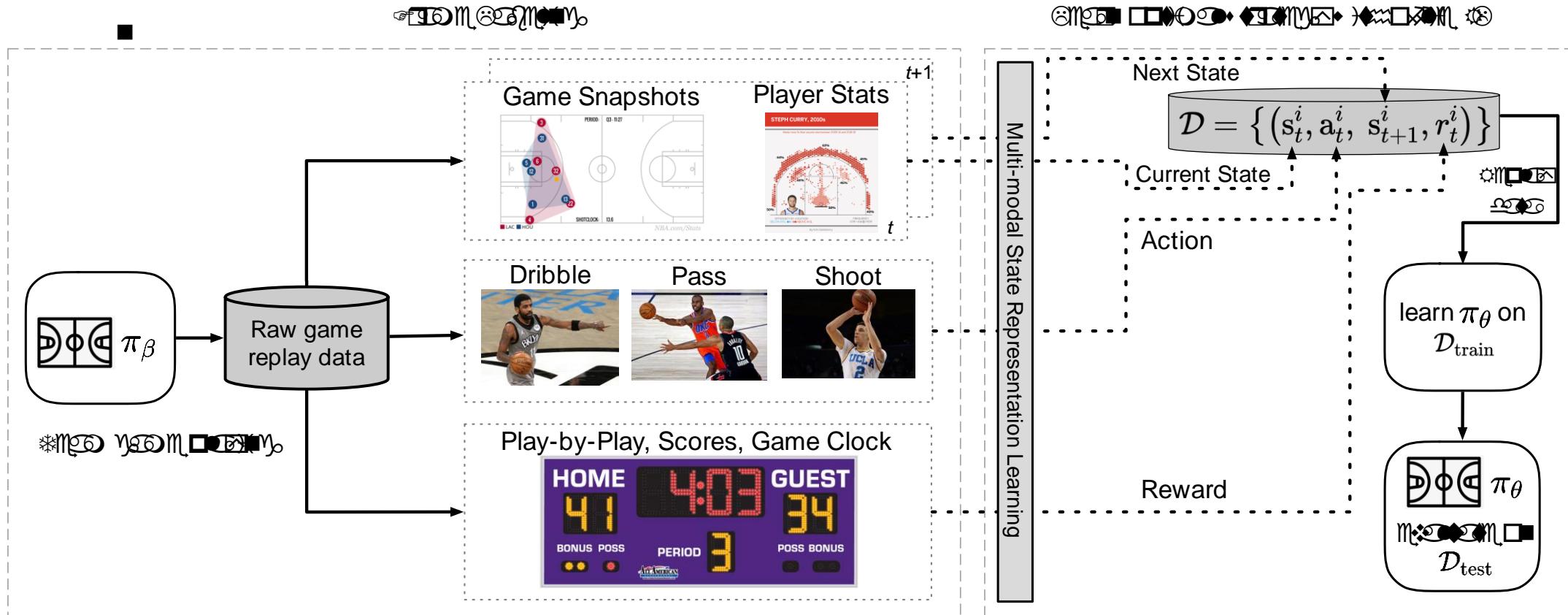
Given a collection of game logs  $\mathcal{D}_{\text{raw}} = \mathcal{D}_{\text{move}} \cup \mathcal{D}_{\text{ppb}} \cup \mathcal{D}_{\text{stat}}$  and an action set  $\mathcal{A}$ , where each  $a \in \mathcal{A}$  is well defined by the discriminative rules on  $\mathcal{D}_{\text{raw}}$ . The task is to assign an appropriate action label  $a$  to every frame in  $\mathcal{D}_{\text{move}}$ . In other words, we aim at producing a policy  $\pi(a | o)$  to output the best action based on the observation related to each frame in  $\mathcal{D}_{\text{move}}$ .

**Reinforcement learning without exploration --> Offline reinforcement learning!**

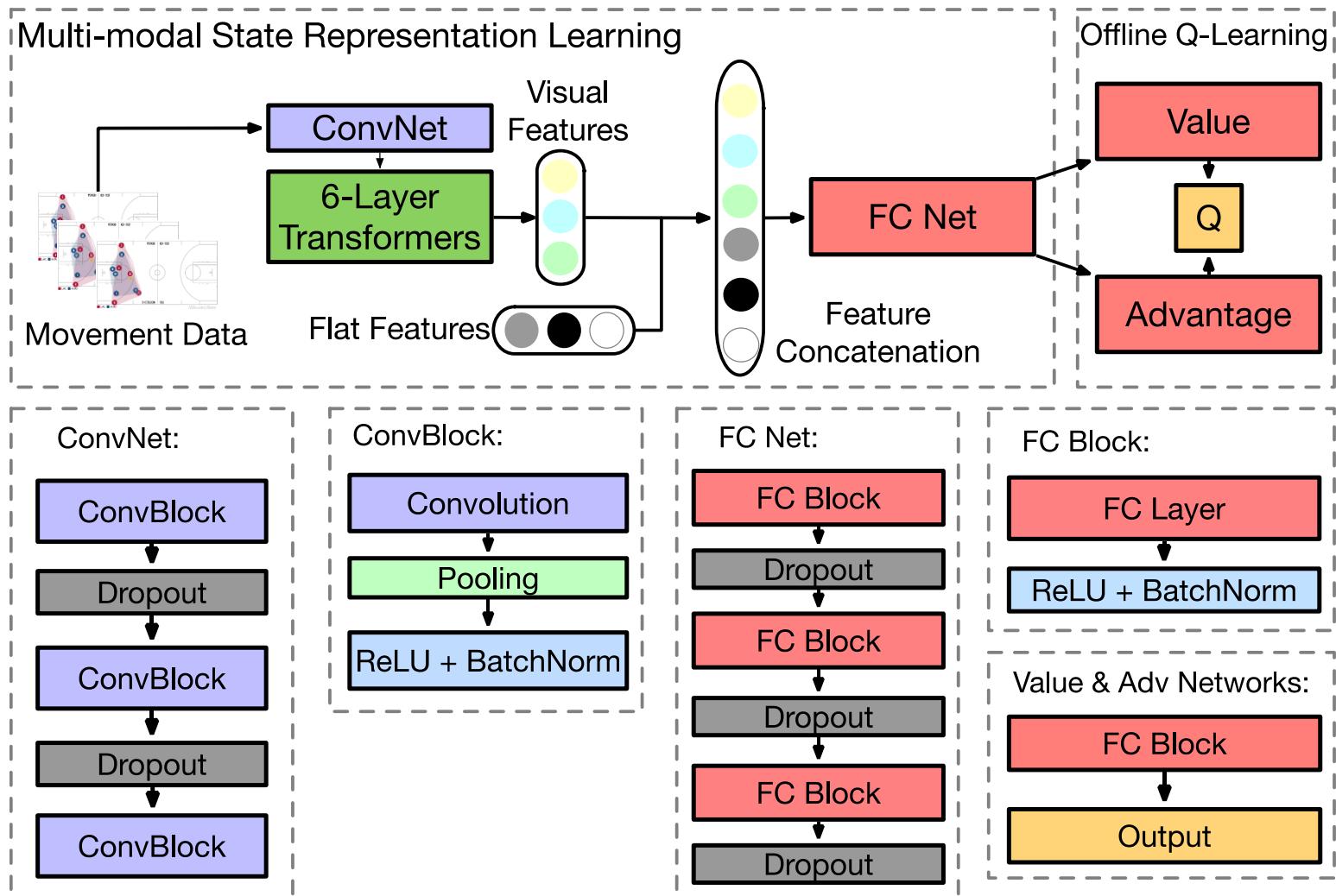
## Why offline RL is challenging?

- No exploration
- The potential cumulative reward is hard to estimate
- Evaluation is hard

# Pipeline Overview



# Architecture Overview



## Reward Function

- Total points scored in this possession
- Shot clock
- Score difference

$$Reward = score + (24.0 - shot\_clock)/24.0 + game\_clock * NB(5, 2/3)$$

## Experimental Settings

- Action Copy
- IS (Importance sampling) – based off-policy evaluation

## Action Copy – binary decision on 3-point attempts

Model	F1 score
Logistic Regression	56.28%
CNN	67.10%
LSTM	68.32%
GRU	67.94%
Transformer	70.43%
Policy Gradient	75.27%
POMDP + Policy Gradient	78.17%
RELIABLE - DQN	76.24%
RELIABLE - POMDP + DQN	<b>81.01%</b>

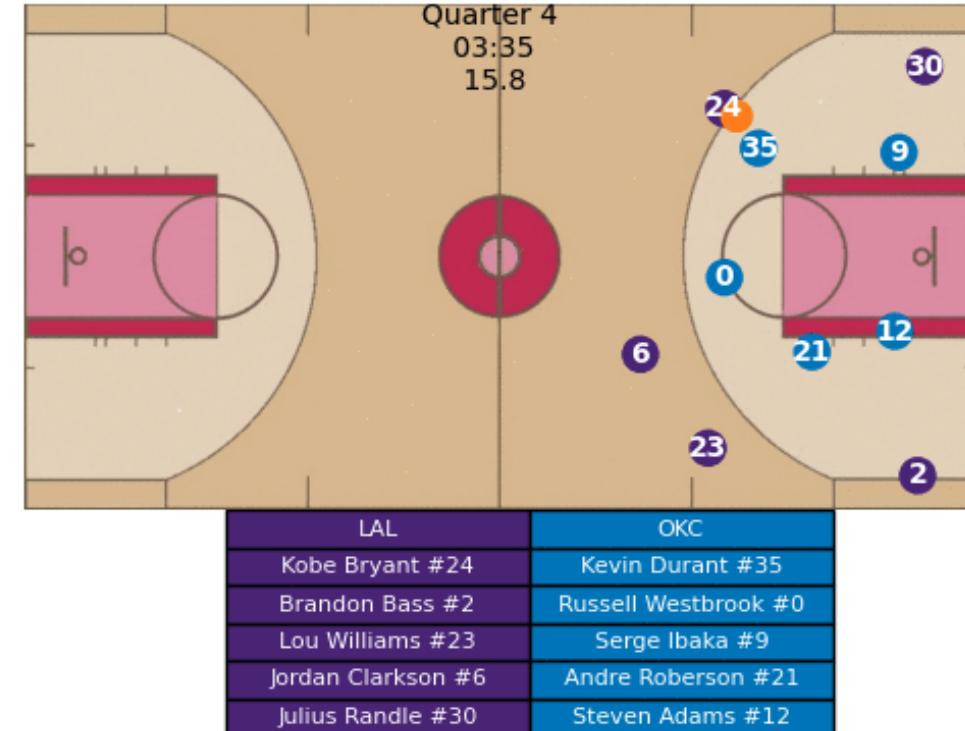
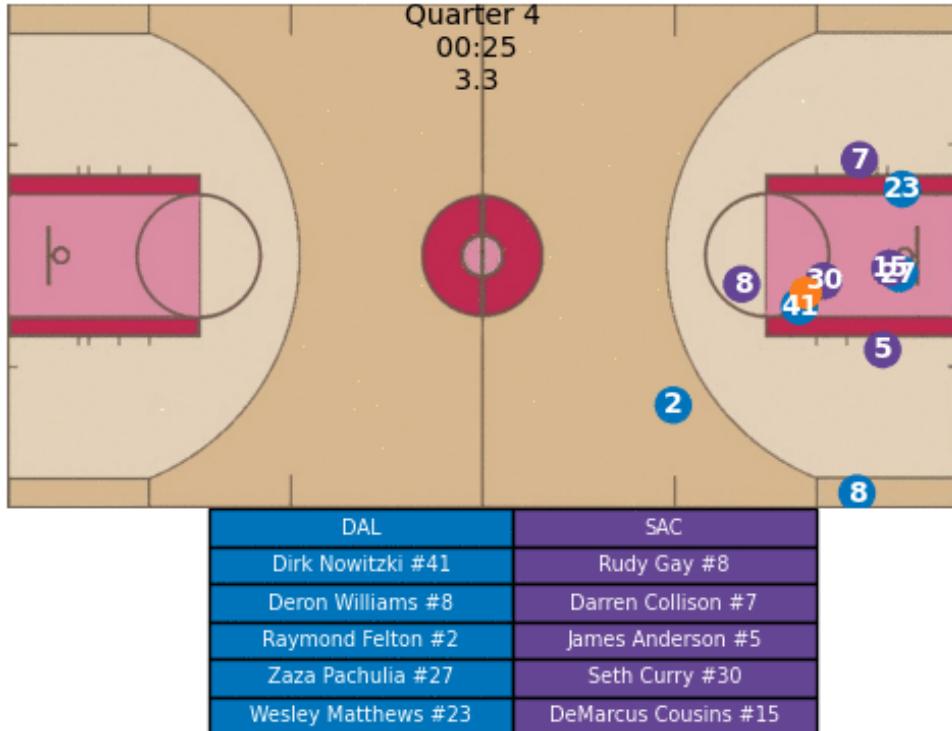
## Action Copy – multi-class decision on {Dribble, Pass, Shoot}

Model	Micro F1	Macro F1
Logistic Regression	35.17%	27.32%
CNN	42.89%	34.71%
LSTM	45.22%	34.75%
GRU	45.74%	34.14%
Transformer	51.20%	37.48%
Policy Gradient	57.36%	40.43%
POMDP + Policy Gradient	70.27%	64.81%
RELIABLE - DQN	60.24%	44.09%
RELIABLE - POMDP + DQN	<b>72.95%</b>	<b>66.90%</b>

## IS (Importance sampling) – based off-policy evaluation

Policy Gradient	81.36
RELIABLE - DQN	94.89
POMDP + Policy Gradient (LSTM)	98.42
RELIABLE - POMDP + DQN (LSTM)	100.28
Season average	102.7
POMDP + Policy Gradient (Transformer)	105.75
RELIABLE - POMDP + DQN (Transformer)	<b>108.16</b>

# Case Studies



## Conclusions

- We propose to formulate the tactical strategy learning of basketball games as solving POMDPs.
- We propose the framework, ReLiable, to apply offline reinforcement learning techniques to solve the POMDP.
- We conduct extensive experiments to showcase that ReLiable can effectively learn good decisions out of replay data without interacting with a real environment.