# Towards Safe and Aligned Embodied AI
# in the Era of Robotics Foundation Models

Thomas (Ran) Tian
rantian@berkeley.edu

🎓 Final-year Ph.D. student at UC Berkeley

Supervised by Prof. Masayoshi Tomizuka and Prof. Andrea Bajcsy

# My Research Trajectory



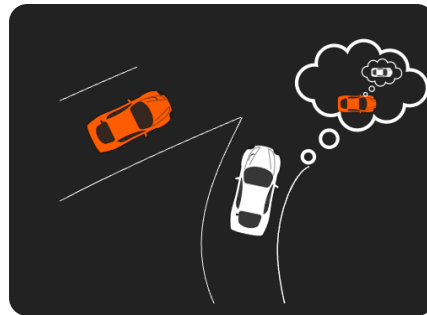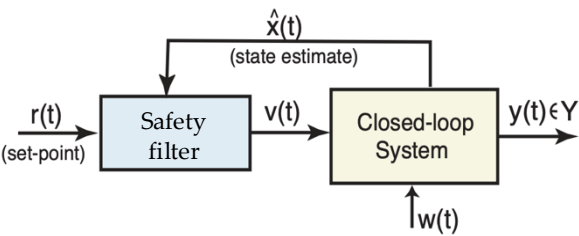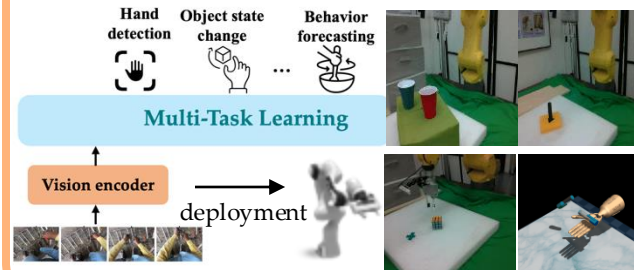**Constrained stochastic control**

**Motion planning**

HRI — Honda Research Institute USA · Qualcomm · WeRide

WAYMO — **Foundation driving model**

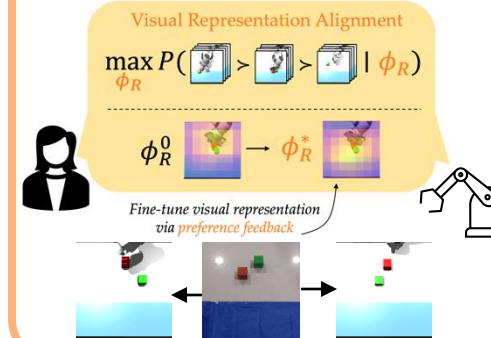Pre-training    Post-training    Efficient deployment
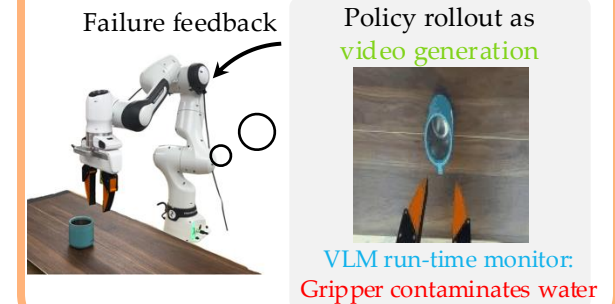
**Dynamical game**

**Safety and alignment**

r(t) (set-point) → Safety filter → v(t) → Closed-loop System → y(t)∈Y

$\hat{x}(t)$ (state estimate)

w(t)

**Representation & policy pre-training**

Hand detection · Object state change · Behavior forecasting · ...

**Multi-Task Learning**

Vision encoder → deployment

**Post-training preference alignment**

Visual Representation Alignment

$\max_{\phi_R} P(\; \cdot\; >\; \cdot\; >\; \cdot\; |\; \phi_R)$

$\phi_R^0 \rightarrow \phi_R^*$

Fine-tune visual representation via preference feedback

**Failure prediction and mitigation beyond collision**

Failure feedback

Policy rollout as video generation

VLM run-time monitor: Gripper contaminates water

# A modern control view of robot autonomy

$$\max_{\pi} \mathbb{E}\left[\sum_{\tau=0}^{T-1} r_R(s_{t+\tau}, \pi(s_{t+\tau})) + V(s_{t+T})\right] \quad \text{s.t.}$$

State space representation of the environment

Task reward

Terminal reward

Dynamics

$$P(s_{t+1}|s_t, a_t)$$

$$P(s_{t+\tau} \in \text{safe set}) \geq \Delta$$

Desired safety guarantee

Model-predictive control

Dynamic programing

Reach-avoid game

Belief space planning
…

# A modern control view of robot autonomy

$$\max_\pi \mathbb{E}\left[\sum_{\tau=0}^{T-1} r_R(s_{t+\tau}, \pi(s_{t+\tau})) + V(s_{t+T})\right] \quad \text{s.t.} \quad \begin{array}{l} P(s_{t+1}|s_t, a_t) \\ P(s_{t+\tau} \in \text{safe set}) \geq \Delta \end{array}$$

**When this works great?**

**Constrained, parsed, and well modeled environment**

Other "agents" are well defined

State is clearly defined

Can know how state evolves given actions of all "agents"

**Have the tools to efficiently solve this problem**

# A modern control view of robot autonomy

$$\max_\pi \mathbb{E}\left[\sum_{\tau=0}^{T-1} r_R(s_{t+\tau}, \pi(s_{t+\tau})) + V(s_{t+T})\right] \quad \text{s.t.} \quad \begin{aligned} &P(s_{t+1}|s_t, a_t) \\ &P(s_{t+\tau} \in \text{safe set}) \geq \Delta \end{aligned}$$

**When this works great?**

**Constrained, parsed, and well modeled environment**

Other "agents" are well defined

State is clearly defined

Can know how state evolves given actions of all "agents"

**Have the tools to efficiently solve this problem**

**But in real world, …**

# A modern control view of robot autonomy

$$\max_\pi \mathbb{E}\left[\sum_{\tau=0}^{T-1} r_R(s_{t+\tau}, \pi(s_{t+\tau})) + V(s_{t+T})\right] \quad \text{s.t.} \quad \begin{array}{l} P(s_{t+1}|s_t, a_t) \\ P(s_{t+\tau} \in \text{safe set}) \geq \Delta \end{array}$$

**When this works great?**

**Constrained, parsed, and well modeled environment**

Other "agents" are well defined

State is clearly defined

Can know how state evolves given actions of all "agents"

**Have the tools to efficiently solve this problem**

**But in real world, ...**

# A modern control view of robot autonomy

$$\max_\pi \mathbb{E}\left[\sum_{\tau=0}^{T-1} r_R(s_{t+\tau}, \pi(s_{t+\tau})) + V(s_{t+T})\right] \quad \text{s.t.} \quad \begin{array}{l} P(s_{t+1}|s_t, a_t) \\ P(s_{t+\tau} \in \text{safe set}) \geq \Delta \end{array}$$

**When this works great?**

**Constrained, parsed, and well modeled environment**

Other "agents" are well defined

State is clearly defined

Can know how state evolves given actions of all "agents"

**Have the tools to efficiently solve this problem**

**But in real world, …**



autonomous, 2x speed

π

# From modern control theory to autonomy 2.0 and beyond

**Beginning:**
**X-MPC with symbolic representations defined by engineers**

State-space representation $\rightarrow$ Model-based approaches $\rightarrow$



Perception models $\rightarrow$

**Structured information**
- o Agent property
  - o Bounding box
  - o Label
  - o Obs status
  - o Risk level
- o Lane property
...

Formulate planning optimization $\rightarrow$

Model-based approaches

# From modern control theory to autonomy 2.0 and beyond

**Beginning:**
**X-MPC with symbolic representations defined by engineers**

State-space representation → Model-based approaches →

**Autonomy 1.0 :**
**[system 1, system 2] with symbolic representations**

State-space representation → Learned motion model / Model-based approaches → supervise/inform/switch

**Human behavior prediction**



**Neural dynamics**
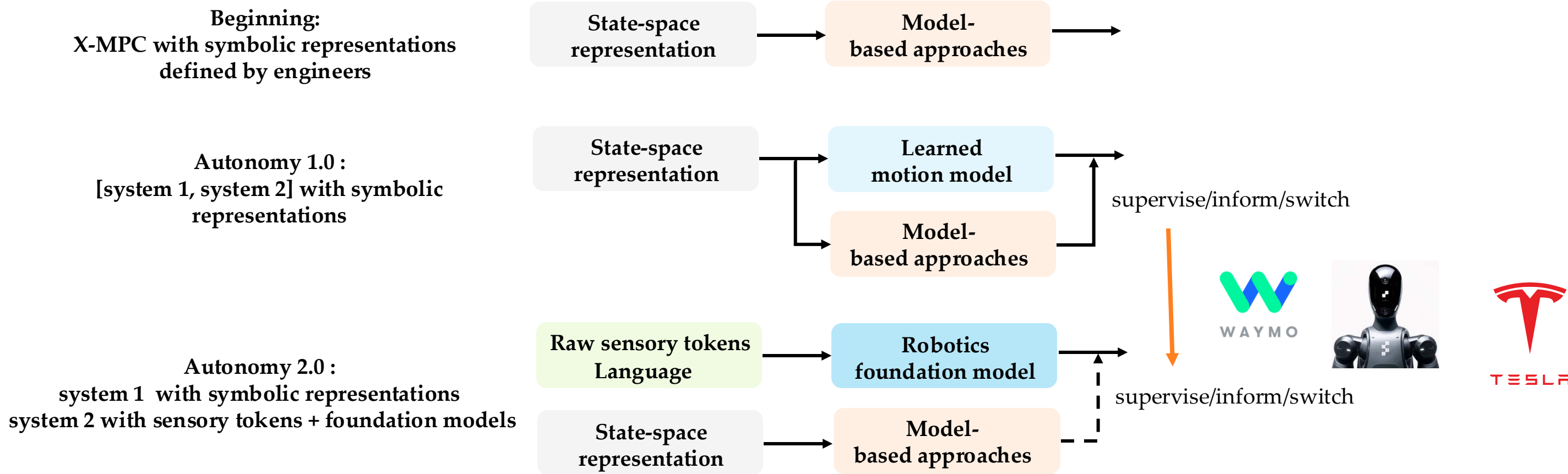


**Control gain scheduling**



· · ·

[Nayakanti, Nigamaa, et al. "Wayformer: Motion forecasting via simple & efficient attention networks." ICRA, 2023.]

[Wang, Changhao, et al. "Safe online gain optimization for cartesian space variable impedance control." CASE. IEEE, 2022.]
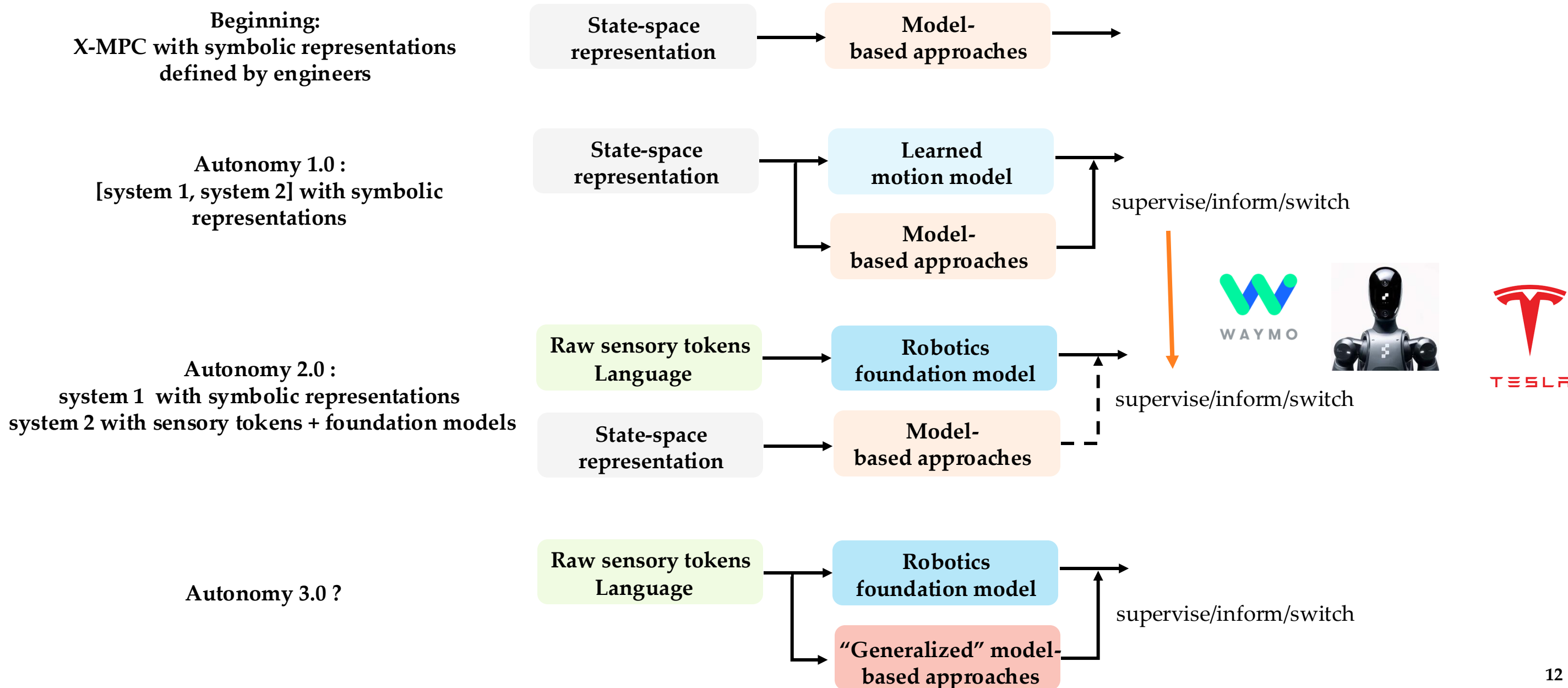
9

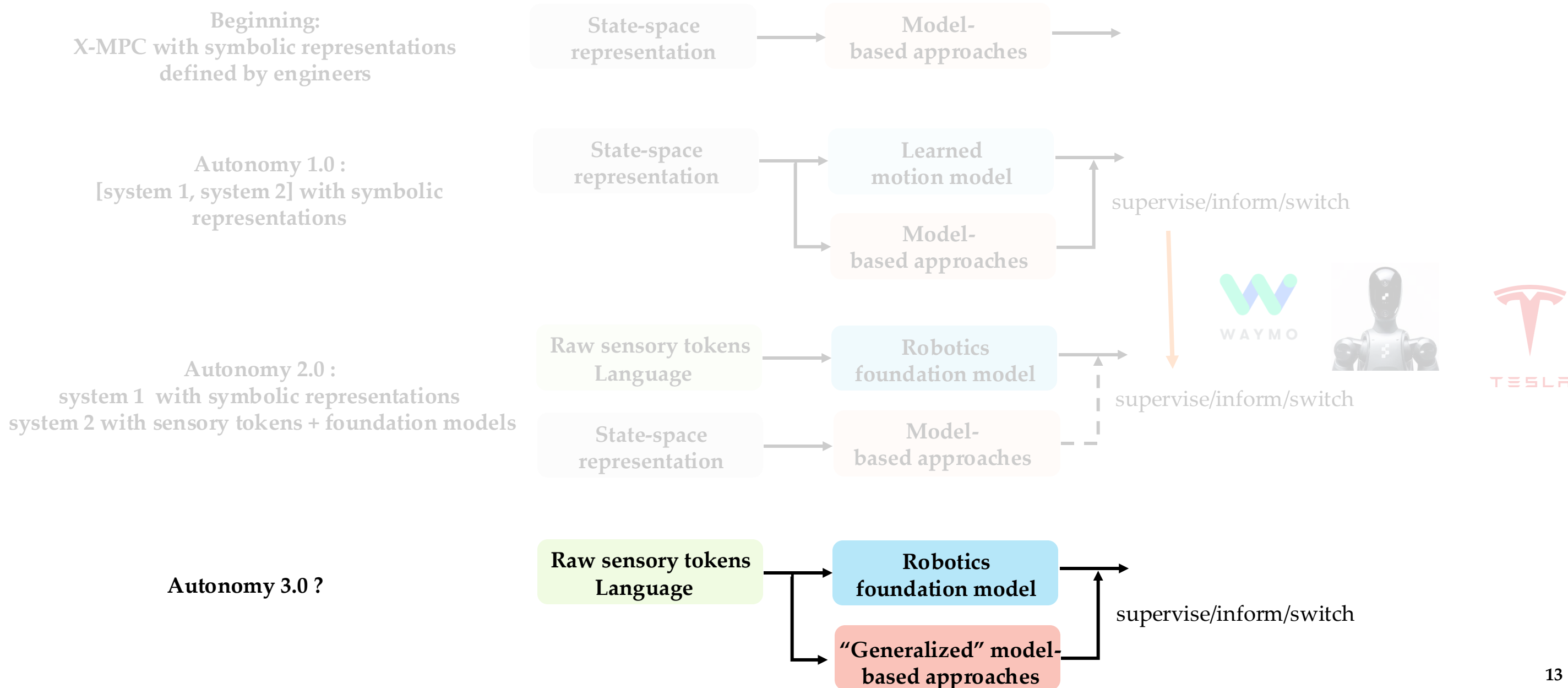# From modern control theory to autonomy 2.0 and beyond

**Beginning:**
**X-MPC with symbolic representations defined by engineers**

State-space representation → Model-based approaches →

**Autonomy 1.0 :**
**[system 1, system 2] with symbolic representations**

State-space representation → Learned motion model

Model-based approaches

supervise/inform/switch

**Autonomy 2.0 :**
**system 1 with symbolic representations**
**system 2 with sensory tokens + foundation models**

Raw sensory tokens Language → Robotics foundation model

State-space representation → Model-based approaches

supervise/inform/switch

# From modern control theory to autonomy 2.0 and beyond

**Beginning:**
**X-MPC with symbolic representations defined by engineers**

State-space representation → Model-based approaches →

**Autonomy 1.0 :**
**[system 1, system 2] with symbolic representations**

State-space representation → Learned motion model / Model-based approaches

supervise/inform/switch

**Autonomy 2.0 :**
**system 1 with symbolic representations**
**system 2 with sensory tokens + foundation models**

Raw sensory tokens Language → Robotics foundation model

State-space representation → Model-based approaches

supervise/inform/switch

WAYMO

TESLA

# From modern control theory to autonomy 2.0 and beyond

**Beginning:**
**X-MPC with symbolic representations defined by engineers**

State-space representation → Model-based approaches →

**Autonomy 1.0 :**
**[system 1, system 2] with symbolic representations**

State-space representation → Learned motion model / Model-based approaches → supervise/inform/switch

**Autonomy 2.0 :**
**system 1 with symbolic representations**
**system 2 with sensory tokens + foundation models**

Raw sensory tokens Language → Robotics foundation model

State-space representation → Model-based approaches

supervise/inform/switch

**Autonomy 3.0 ?**

Raw sensory tokens Language → Robotics foundation model / "Generalized" model-based approaches → supervise/inform/switch

12

# From modern control theory to autonomy 2.0 and beyond

**Beginning:**
**X-MPC with symbolic representations defined by engineers**

State-space representation → Model-based approaches →

**Autonomy 1.0 :**
**[system 1, system 2] with symbolic representations**

State-space representation → Learned motion model / Model-based approaches

supervise/inform/switch

**Autonomy 2.0 :**
**system 1  with symbolic representations**
**system 2 with sensory tokens + foundation models**

Raw sensory tokens Language → Robotics foundation model

State-space representation → Model-based approaches

supervise/inform/switch

**Autonomy 3.0 ?**

Raw sensory tokens Language → Robotics foundation model / "Generalized" model-based approaches

supervise/inform/switch

13

# Behavior cloning for robot learning



[Toyota Research Institute unveils breakthrough in Teaching Robots New Behaviors, https://www.tri.global/news/toyota-research-institute-unveils-breakthrough-teaching-robots-new-behaviors ]

# Behavior cloning for robot learning

**Let's remember what the expert did and copy them!**

$$\max_{\theta} \mathbb{E}\left[\mathbb{P}_{\theta}(\boldsymbol{a}^*_{0:T}|\boldsymbol{o}_0; \text{context})\right]$$

$$\{(\text{context}, \boldsymbol{o}_0, \boldsymbol{a}^*_{0:T})^1, ....\}$$

# 10 years ago – CNN based Policy Model



[Bojarski, Mariusz, et al. "End to end learning for self-driving cars." *arXiv preprint arXiv:1604.07316* (2016)]

# 2 years ago – Diffusion Policy Model



[Chi, Cheng, et al. "Diffusion policy: Visuomotor policy learning via action diffusion." The International Journal of Robotics Research (2023): 0278364241273668.]

# Large-scale data is a key factor for robotics foundation models



$$\{(\text{observation}, \text{action}, \text{task spec})_t \}$$

[Open-X collaboration, **Ran Tian**, et al. , Open X-Embodiment: Robotic Learning Datasets and RT-X Models, ICRA'24, best-paper award]

# Now – Vision-Language-Action Robotics Foundation Model

*Q: Wha should the robot do to pick up the chip bag?*

**Large Language Model**

*A: The robot should use its gripper to pick up the bag*

[Brohan, Anthony, et al. "Rt-2: Vision-language-action models transfer web knowledge to robotic control." arXiv preprint arXiv:2307.15818 (2023)]

# Now – Vision-Language-Action Robotics Foundation Model

Q: *Wha should the robot do to pick up the chip bag?*

**Vision-Language-Action Model**

A: *The robot should use its gripper to pick up the bag*

A: = 132 114 128 5 25 156

De-tokenize

Δ T = [0.1, -0.2, 0]
Δ R = [10˚, 25˚, -7˚]

Robot action

Put strawberry into the correct bowl (2x speed)

[Brohan, Anthony, et al. "Rt-2: Vision-language-action models transfer web knowledge to robotic control." arXiv preprint arXiv:2307.15818 (2023)]

# Now – Vision-Language-Action Robotics Foundation Model



Q: *Wha should the autonomous vehicle do?*

<Multi-view video, HD-map, ..., Symbolic representation>

Object-centric Tokenization

Scene-centric Tokenization

**Vision-Language-Action Model**

A: *The robot should use its gripper to pick up the bag*

A: = 132 114 128 5 25 156

De-tokenize

Δ T = [0.1, -0.2, 0]
Δ R = [10˚, 25˚, -7˚]

Robot action

[**Ran Tian**, et al. "Tokenize the world into object-level knowledge to address long-tail events in autonomous driving." CoRL, 2024.]

# Shortest Path to Generalist Robots?

**Scaling law seems to work in robotics**



**Commercialization of large-scale robot data collection**





**Open-sourced model / training pipeline**



**More affordable "one-click" fine-tunning API**

# Imitation is a "proxy" of the true training objective

**Let's remember what the expert did and copy them!**

$$\max_{\theta} \mathbb{E}\left[\mathbb{P}_{\theta}(\boldsymbol{a}_{0:T}^{*}|\boldsymbol{o}_0; \text{context})\right]$$

$\neq$

*Safety > comfort, progress, etc*

## Miss-Alignment

By optimizing an **incomplete** or **mis-specified** objective, these models lead to undesirable behaviors at best and safety hazards at worst!

Hand me a bag of chips

*I don't want my chips crashed...*

$$\max_{\theta} \mathbb{E}\left[\mathbb{P}_{\theta}(\boldsymbol{a}_{0:T}^{*}|\boldsymbol{o}_{0}; \text{context})\right]$$

# Reinforcement Learning from Human Feedback for Post-training Preference Alignment

**Step 2: Sample generations from the model**

(A) (B) (C)

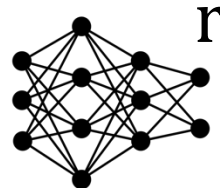**Step 5: Train the model to max. the reward**

max r

**Model-based principles come back!**

**Generative Model**

**Step 1: Train a generative model via BC**

(A) > (B) > (C)

**Step 3: Ask humans to grade**

r

**Step 4: Learn a reward function**

[Christiano, Paul F., et al. "Deep reinforcement learning from human preferences." Advances in neural information processing systems 30 (2017).]

# Reinforcement Learning from Human Feedback for Post-training Preference Alignment

Predominant alignment mechanism in *non-embodied* domains



**Text-image gen. :** A cake in the city.

Pre-trained    After alignment



**Text-video gen. :** An animated sneaker is playing basketball.

Pre-trained    After alignment

RLHF has yet to achieve the same impact in aligning robotics motion generation models

# Reinforcement Learning from Human Feedback
# for Post-training Preference Alignment



**Step 2: Sample generations from the model**

$(A) \quad (B) \quad (C)$

$(A) > (B) > (C)$

**Step 3: Ask humans to grade**

**Step 5: Train the model to max. the reward**

$\max r_R$

**Motion Generation Model**

**Step 1: Train a generative model via BC**

$r$

**Step 4: Learn a reward function**

# Reinforcement Learning from Human Feedback
# for Post-training Preference Alignment

**Step 2: Sample generations from the model**



Ⓐ    Ⓑ    Ⓒ

**Step 5: Train the model to max. the reward**

$$\max r_R(\text{🦾}, \text{task})$$

**Motion Generation Model**

**Step 1: Train a generative model via BC**

Ⓐ > Ⓑ > Ⓒ

**Step 3: Ask humans to grade**

$r$

**Step 4: Learn a reward function**

Learning a high-quality visual reward function requires an impractically large amount of human preference feedback

# Reinforcement Learning from Human Feedback for Post-training Preference Alignment



**Step 2: Sample generations from the model**

**Step 5: Train the model to max. the reward**

$$\max r_R(\succ, task)$$

**Motion Generation Model**

**Step 1: Train a generative model via BC**

(A) > (B) > (C)

**Step 3: Ask humans to grade**

$r$

**Step 4: Learn a reward function**

Goal: Maximizing Alignment with Minimal Feedback!

$r_R(\ )$

Policy Behavior $o^\pi$

" similar? "

Encode robot
video with $\phi_H$

Human expert

Preferred Behavior $o^+$

33

$r_R(\quad)$

$\rho(\phi_H(\boldsymbol{o}^\pi))$

*Policy Feature Distribution*

$r_H(\phi_H(\quad))$

*Encode robot video with $\phi_H$*

*match* $\begin{cases} \text{Forward KL divergence} \\ \text{Reverse KL divergence} \\ \text{JS divergence} \\ \vdots \end{cases}$

$\rho(\phi_H(\boldsymbol{o}^+))$

*Preferred Feature Distribution*

*Visual Representation Alignment*

$$\underset{\phi_R}{\text{argmin}}\ \text{diff}(\ \phi_R, \phi_H)$$

$\text{r}_R(\phi_R(\ ))$

$\rho(\phi_H(\boldsymbol{o}^\pi))$

*Policy Feature Distribution*

$\text{r}_H(\phi_H(\ ))$

*Encode robot video with $\phi_H$*

*match*

$\rho(\phi_H(\boldsymbol{o}^+))$

*Preferred Feature Distribution*

35

$$\max r_R(\; ,\text{task})$$

**Motion Generation Model**

$$A > B > C$$

r

Learning reward end-to-end

**Key idea:** Allocate human budget *exclusively* to *align* visual representations

Visual Representation Alignment

# Formalizing Visual Representation Alignment

*How to formally describe and compare two encoders?*



$$\underset{\phi_R}{\mathrm{argmin}}\ \mathrm{diff}(\ \phi_R, \phi_H)$$

# Formalizing Visual Representation Alignment

*How to formally describe and compare two encoders?*

Triplet-based *representation space* [Sucholutsky & Griffiths, NeurIPS 2023]



$$\phi \longrightarrow S_\phi := \{ ( \mathbf{o}^i, \mathbf{o}^j, \mathbf{o}^k ) : d\left( \phi(\mathbf{o}^i), \phi(\mathbf{o}^j) \right) < d\left( \phi(\mathbf{o}^i), \phi(\mathbf{o}^k) \right), \mathbf{o}^{i,j,k} \in \Xi \}$$

$\mathbf{o}^i$   $\mathbf{o}^j$   $\mathbf{o}^k$

anchor   positive   negative

**$\mathbf{o}^i$ is closer to $\mathbf{o}^j$ than to $\mathbf{o}^k$**

*learning $\phi_R$ which minimizes the difference between two agents' representation spaces (S)*

$$\underset{\phi_R}{\text{argmin}}\ \text{dis}(S_{\phi_R}, S_{\phi_H})$$

$\phi_R$   $z$

**1) Don't have direct access to $\phi_H$**
**2) $S_{\phi_H}$ is extremely large in space of videos**

$$\min_{\phi_R} \ell(S_{\phi_R}, S_{\phi_H})$$

$$\min_{\phi_R} \ell(S_{\phi_R}, \tilde{S}_{\phi_H})$$

Query the end user

Preference Feedback



*I don't want my chips crushed!*

*End-User*

$$\min_{\phi_R} \ell(S_{\phi_R}, \tilde{S}_{\phi_H})$$

Preference Feedback

I don't want my chips crushed!

End-User

$\tilde{S}_{\phi_H}$

Visual Representation Alignment

better    worse    most preferred

$$\max_{\phi_R} \mathbb{P}(\quad \succ \quad | \quad \phi_R)$$

$\mathbf{o}^j \qquad \mathbf{o}^k \qquad \mathbf{o}^i$

$\phi_R \to z$

Pre-trained Vision Encoder

Bradley-Terry model

$$\frac{e^{-d(\phi_R(\mathbf{o}^i), \phi_R(\mathbf{o}^j))}}{e^{-d(\phi_R(\mathbf{o}^i), \phi_R(\mathbf{o}^j))} + e^{-d(\phi_R(\mathbf{o}^i), \phi_R(\mathbf{o}^k))}}$$

Two equally **preferred** behaviors should have **similar** feature representations.

42

Preference Feedback

*I don't want my chips crushed!*

*End-User*

$\tilde{S}_{\phi_H}$

Visual Representation Alignment

better · worse · most preferred

$$\max_{\phi_R} \mathbb{P}(\ \mathbf{o}^j \succ \mathbf{o}^k \mid \mathbf{o}^i \quad \phi_R\ )$$

$\phi_R \rightarrow z$

*Pre-trained Vision Encoder*

$\phi_R^*$

Generative Model
Preference Alignment

$\phi_R^* \quad z \quad \pi_R \rightarrow a_R$

$$\max_{\pi} \quad r(\boldsymbol{o}^\pi \mid \phi_R^*)$$

43

# Representation-Aligned Preference-Based Learning (RAPL)

[Tian et al. "What Matters to You? Towards Visual Representation Alignment for Robot Learning". ICLR 2024.]

[Tian et al. "Maximizing Alignment with Minimal Feedback: Efficiently Learning Rewards for Visuomotor Robot Policy Alignment". arxiv 2025.]

# On the Value of Aligned Representation in Visual Reward Learning

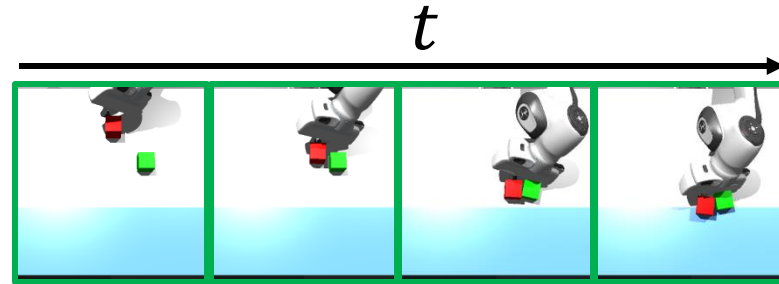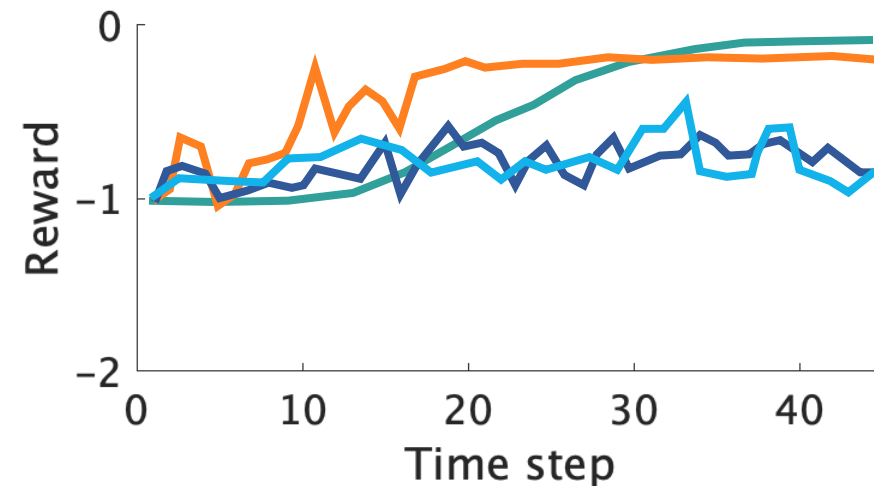$r^*(\phi_H(o))$: push objects move **efficiently** to **goal region**

$t$



## Reward of *Good* Behavior



$r^*$ Ground truth reward

Ours (rep. alignment first then reward pred.)

# On the Value of Aligned Representation in Visual Reward Learning

$r^*(\phi_H(o))$: push objects move **efficiently** to **goal region**



$t$

Reward of *Good* Behavior



- $r^*$ Ground truth reward
- Ours (rep. alignment first then reward pred.)
- MVP representation [1]

[1] He et al., "Masked autoencoders are scalable vision learners." CVPR 2022

# On the Value of Aligned Representation in Visual Reward Learning

$r^*(\phi_H(o))$:  push objects move **efficiently** to **goal region**

$t$



## Reward of *Good* Behavior



$r^*$ Ground truth reward

Ours (rep. alignment first then reward pred.)

MVP representation

Dino representation [1]

[1] Oquab et al., "DINOv2: Learning Robust Visual Features without Supervision."

# On the Value of Aligned Representation in Visual Reward Learning

$r^*(\phi_H(o))$: push objects move **efficiently** to **goal region**

$t$



## Reward of *Good* Behavior



**DINO-WM: World Models on Pre-trained Visual Features enable Zero-shot Planning**

Gaoyue Zhou [1]  Hengkai Pan [1]  Yann LeCun [1,2]  Lerrel Pinto [1]

**Real-World Robot Learning with Masked Visual Pre-training**

Ilija Radosavovic*  Tete Xiao*  Stephen James  Pieter Abbeel  Jitendra Malik[†]  Trevor Darrell[†]

University of California, Berkeley

**Abstract:** In this work, we explore self-supervised visual pre-training on images from diverse, in-the-wild videos for real-world robotic tasks. Like prior work, our visual representations are pre-trained via a masked autoencoder (MAE), frozen, and then passed into a learnable control module. Unlike prior work, we show that the pre-trained representations are effective across a range of real-world robotic tasks and embodiments. We find that our encoder consistently outperforms CLIP (up to 75%), supervised ImageNet pre-training (up to 81%), and training from scratch (up to 81%). Finally, we train a 307M parameter vision transformer on a massive collection of 4.5M images from the Internet and egocentric videos, and demonstrate clearly the benefits of scaling visual pre-training for robot learning.

**Keywords:** Self-supervised Learning, Visual Representations, Robot Learning

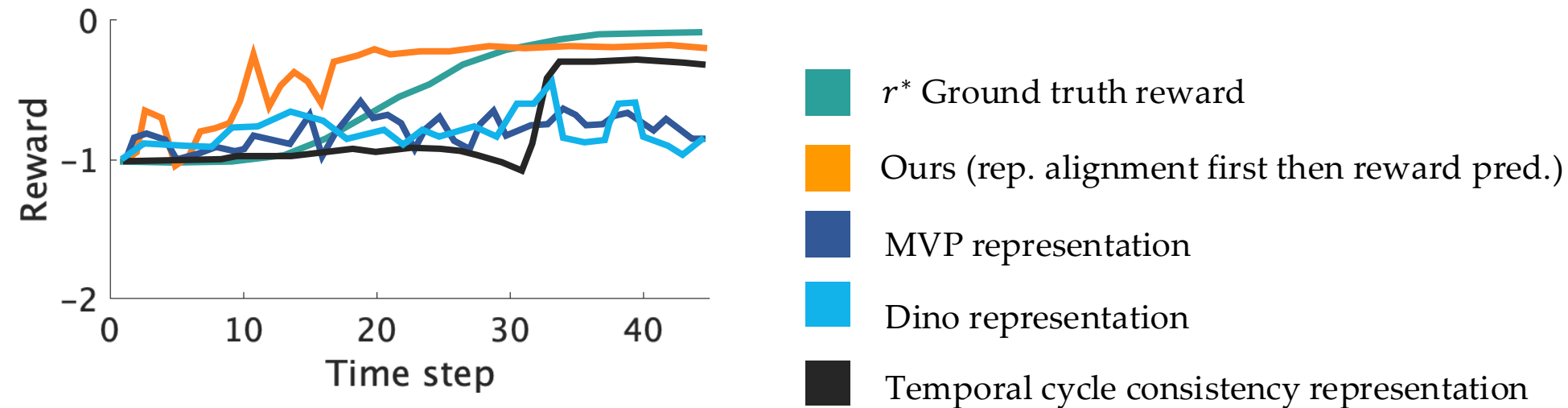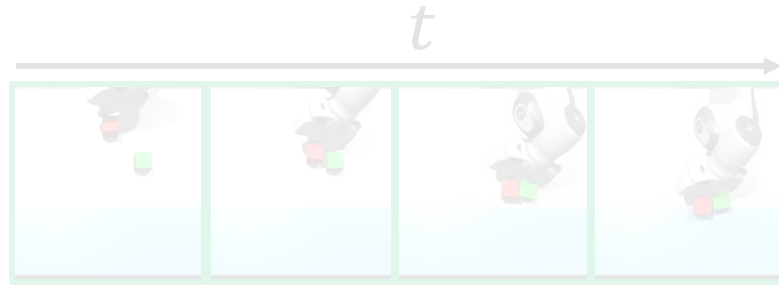Serve as **latent state** representation for planning

[1] Oquab et al., "DINOv2: Learning Robust Visual Features without Supervision."

# On the Value of Aligned Representation in Visual Reward Learning

$r^*(\phi_H(o))$: push objects move **efficiently** to **goal region**
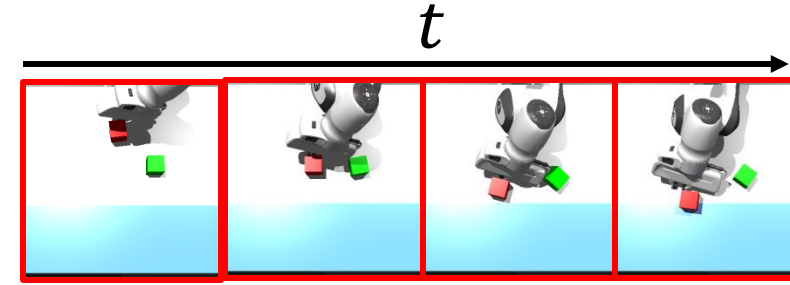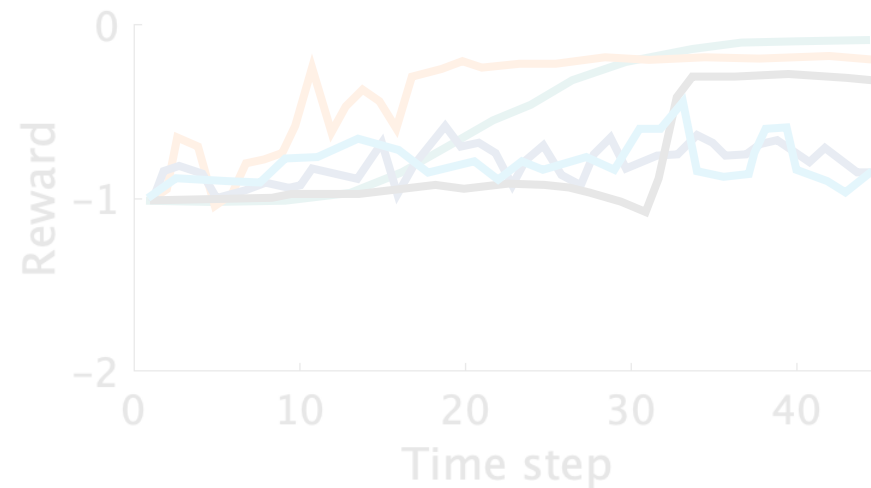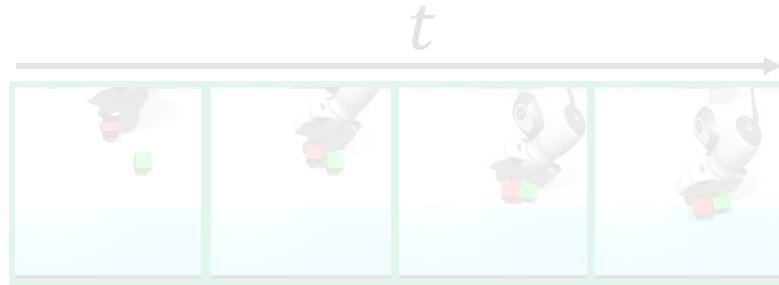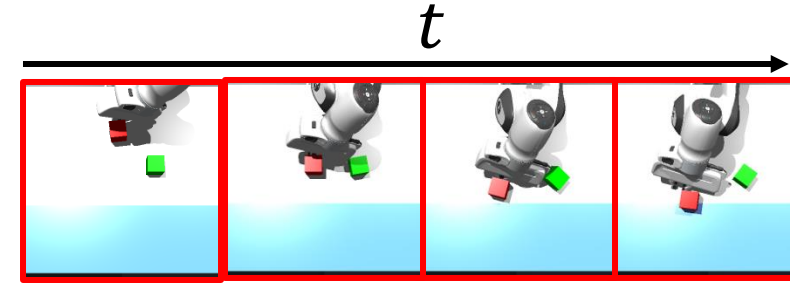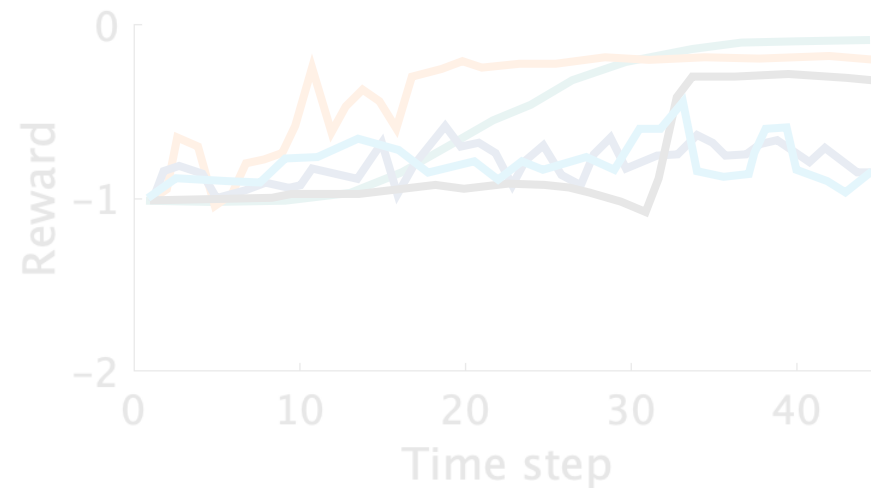


Reward of *Good* Behavior

Reward

Time step

- $r^*$ Ground truth reward
- Ours (rep. alignment first then reward pred.)
- MVP representation
- Dino representation
- Temporal cycle consistency representation

$\phi_{TCC}[1]$

[1] Dwibedi, et al. "Temporal cycle-consistency learning." CVPR 2019.

# On the Value of Aligned Representation in Visual Reward Learning

$r^*(\phi_H(o))$: push objects move **efficiently** to **goal region**
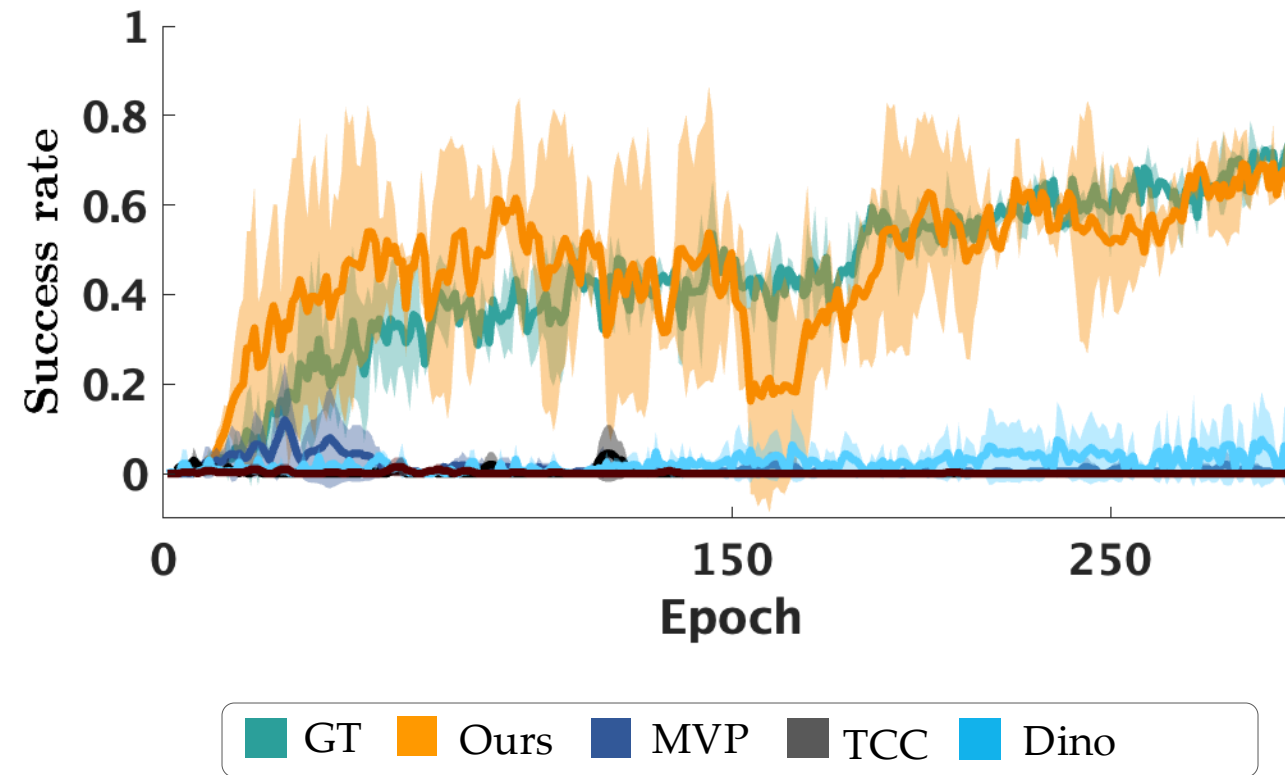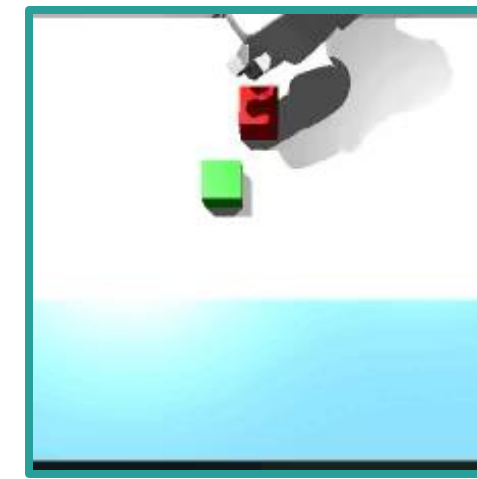


Reward of *Good* Behavior

Reward of *Disliked* Behavior

GT    Ours    MVP    TCC    Dino

# On the Value of Aligned Representation in Visual Reward Learning

$r^*(\phi_H(o))$:  push objects move **efficiently** to **goal region**

# On the Value of Aligned Representation in Visual Reward Learning
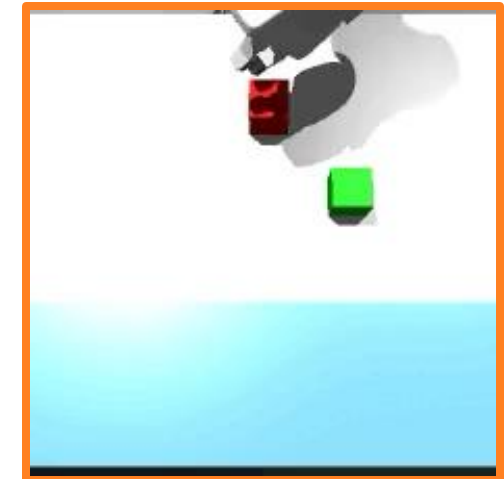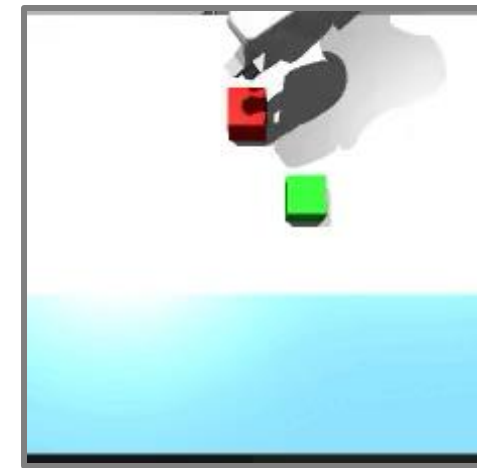


**Insight**
Pre-trained unaligned representations might miss important features that matter to the task
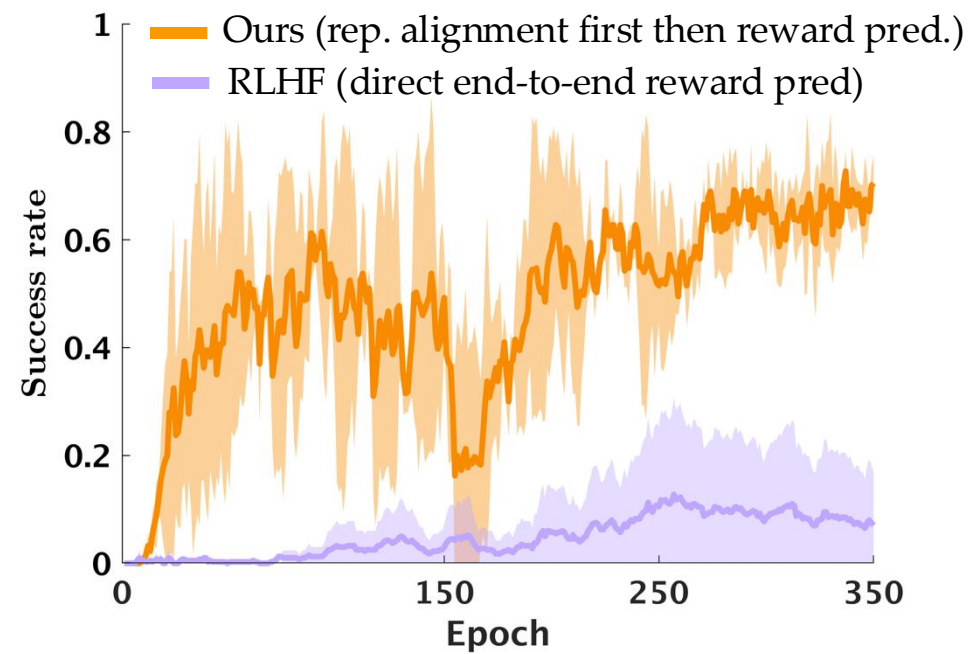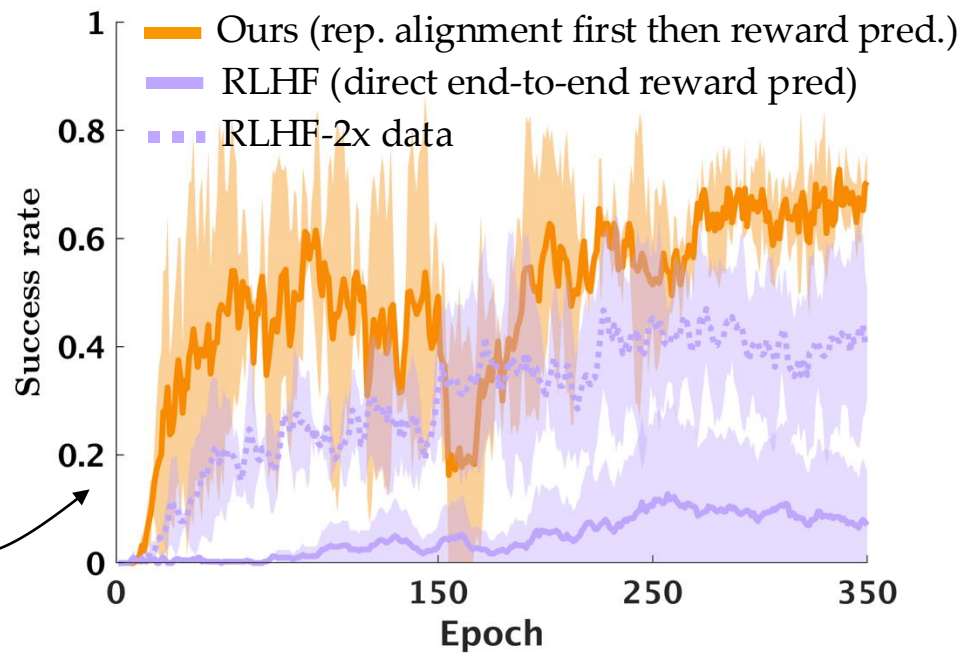


GT

RAPL

TCC

Dino

55

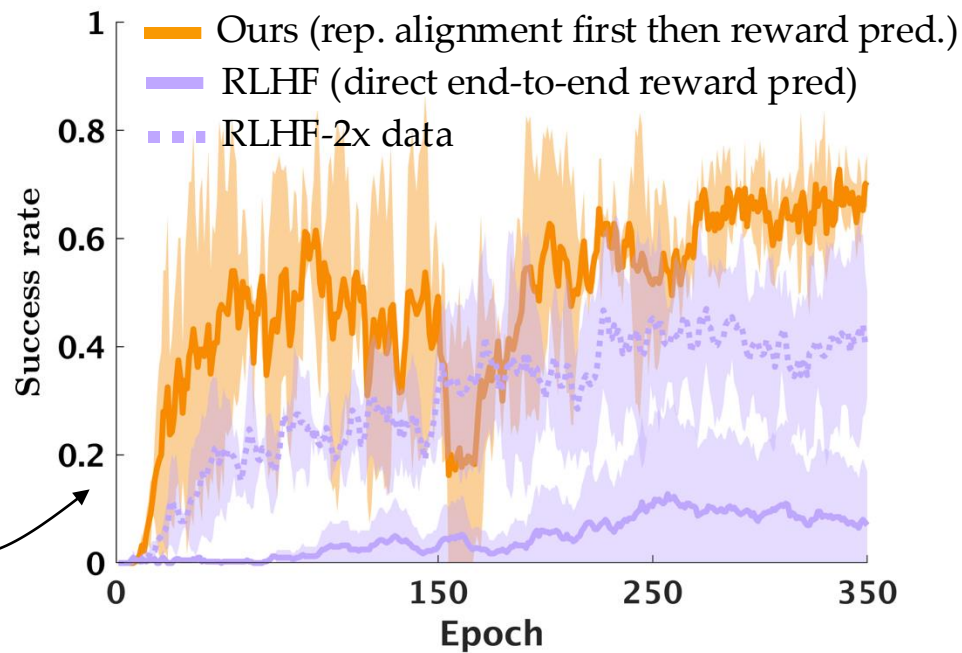# On the Value of Aligned Representation – *Sample Efficiency*

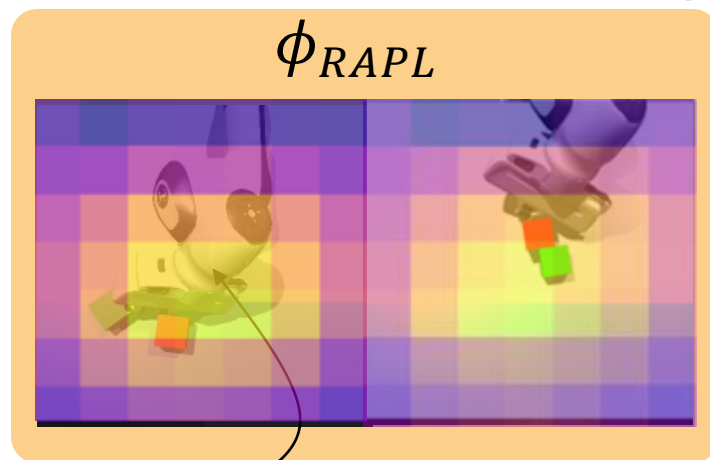# On the Value of Aligned Representation – *Sample Efficiency*



*RAPL outperforms RLHF by 75% with 50% less preference data*

# On the Value of Aligned Representation – *Sample Efficiency*
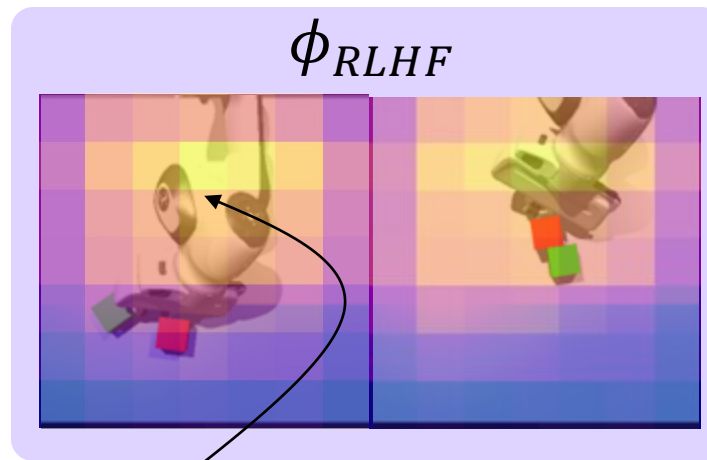


*RAPL outperforms RLHF by 75% with 50% less preference data*

$\phi_{RAPL}$

$\phi_{RLHF}$

Contribution to final embed.

high | low

*Attending to task-relevant goal region & objects*

*Attending to embodiment*

58

# On the Value of Aligned Representation – *Sample Efficiency*



*RAPL outperforms RLHF by 75% with 50% less preference data*

**Insight**

Representation alignment **reduces** the reliance on human preference feedback for achieving effective preference alignment.

$\phi_{RAPL}$

$\phi_{RLHF}$

$\phi_{RLHF}$-2x data

Contribution to final embed.

high          low

*Attending to task-relevant goal region & objects*

*Attending to embodiment*

# Preference Alignment of Real-World Robot Visuomotor Policies

**Behavior Cloning**

$$\max \mathbb{P}(\boldsymbol{a}_t^{\textbf{demo}} | \boldsymbol{o}_t; \text{task})$$

**Pick up chips**

Crush the chips

**Pick up cup**

Contaminate water

**Pick up and place fork**

Drop the fork

# Preference Alignment of Real-World Robot Visuomotor Policies



No access to high-fidelity simulators in deployment setting.

A > B > C

$\underset{\phi_R}{\text{argmin}}\ \text{diff}(\phi_R, \phi_H)$

Visual Representation Alignment

**Motion Generation Model**

Reward via Feature Matching

$\phi_R^*(o^\pi)$

*Policy Feature Distribution*

$\updownarrow$ matching

$\phi_R^*(o^+)$

*Preferred Feature Distribution*

$\max\ r_R(o^{\pi_{ref}}; \phi_R^*)$

# Preference Alignment of Real-World Robot Visuomotor Policies



Need a huge number of rankings!

DPO: update generative model without RL!

**Contrastive Learning**

$$\max_{\pi} \mathbb{P}(a_A^{\pi_{ref}} \succ a_B^{\pi_{ref}})$$

**Motion Generation Model**

A > B > C

[Rafailov, Rafael, et al. "Direct preference optimization: Your language model is secretly a reward model." Advances in Neural Information Processing Systems 36 (2023): 53728-53741.]

# Preference Alignment of Real-World Robot Visuomotor Policies

**Key idea**: leverage the reward built from the aligned representation to scalably generate synthetic rankings.

# Preference Alignment of Real-World Robot Visuomotor Policies

**Imitation Learning**

$$\max \mathbb{P}(\boldsymbol{a}_t^{\mathbf{demo}}|\boldsymbol{o_t}; \text{task})$$

*Before alignment*

$$\max \mathbb{P}(\boldsymbol{a}_+^{\pi_{ref}} \succ \boldsymbol{a}_-^{\pi_{ref}}; \phi_R^*, \text{task})$$

**Preference Alignment**

*After alignment*



**Pick up chips**

Crush the chips

**Pick up cup**

Contaminate water

**Pick up and place fork**

Drop the fork

Hold the packaging by its edges

Pick the cup by the handle

Gently place the fork

64

# Preference Alignment of Real-World Robot Visuomotor Policies

*Result – alignment performance*

**Ours** achieves **2x** (in average) better alignment scores compared to baselines under same amount of human budget



Legend:
- Ref. IL
- Ours (rep. alignment first then reward pred.)
- Direct reward pred. (conventional RLHF)
- MVP
- TCC
- R3M

Alignment score

Cup    Fork    Bag

$$\frac{\# \text{ most-likely mode is preferred } (\textbf{\textit{graded by human}})}{\# \text{ test task configurations}}$$

# Let's Scale it up to Multi-agent!

**Large scale traffic simulation**



agents

Predicted action tokens

**Motion Generation Model**

Past obs. tokens    Past action tokens

# Imitation-learning based Traffic Simulation Model



$\mathbb{P}(\boldsymbol{a}_0; c)$  $\mathbb{P}(\boldsymbol{a}_1; \hat{\boldsymbol{a}}_0, c)$ ... $\mathbb{P}(\boldsymbol{a}_T; \hat{\boldsymbol{a}}_{<T}, c)$

**Traffic Simulation Model**

**start**  $\hat{\boldsymbol{a}}_0$ ... $\hat{\boldsymbol{a}}_{T-1}$

$\begin{bmatrix} \mathbf{a}_0 & \mathbf{a}_1 & \cdots & \mathbf{a}_T \end{bmatrix}$

(action = {long. accel, lat. accel})

# Imitation-learning based Traffic Simulation Model

$\mathbb{P}(\boldsymbol{a}_0; c)$  $\mathbb{P}(\boldsymbol{a}_1; \hat{\boldsymbol{a}}_0, c)$  ...  $\mathbb{P}(\boldsymbol{a}_T; \hat{\boldsymbol{a}}_{<T}, c)$

**Imitation Learning**

$$\max_{\theta} \Pi_t^T \mathbb{P}(\boldsymbol{a}_t | \boldsymbol{a}_{t-1}, ..., \boldsymbol{a}_0; c)$$

**Traffic Simulation Model**

**start**  $\hat{\boldsymbol{a}}_0$  ...  $\hat{\boldsymbol{a}}_{T-1}$



Realism Score

SMART-large

Behavior-GPT

SMART

GUMP

IL pre-trained

0.76

0.74

0.72

0.70

1    5    10    >100

Model Capacity (M)

# Imitation-learning based Traffic Simulation Model



$\mathbb{P}(\boldsymbol{a}_0; c)$    $\mathbb{P}(\boldsymbol{a}_1; \hat{\boldsymbol{a}}_0, c)$   ...   $\mathbb{P}(\boldsymbol{a}_T; \hat{\boldsymbol{a}}_{<T}, c)$

**Imitation Learning**

$$\max_{\theta} \Pi_t^T \mathbb{P}(\boldsymbol{a}_t | \boldsymbol{a}_{t-1}, \dots, \boldsymbol{a}_0; c)$$

**Traffic Simulation Model**

start    $\hat{\boldsymbol{a}}_0$   ...   $\hat{\boldsymbol{a}}_{T-1}$

SMART-large

0.76

Behavior-GPT

SMART

0.74

GUMP

Realism Score

0.72

IL pre-trained

0.70

1    5    10    >100

**Model Capacity (M)**

# Preference-Alignment of Traffic Simulation Model



[Ran Tian et al. "Direct Post-Training Preference Alignment for Multi-Agent Motion Generation Models Using Implicit Feedback from Pre-training Demonstrations." ICLR 2025, Spotlight.]

# Preference-Alignment with *Hand-designed* Features



$\phi_p(\boldsymbol{o}^\pi)$

*Policy Feature Distribution*

$\updownarrow$ matching

$\phi_p(\boldsymbol{o}^+)$

*Preferred Feature Distribution*

$$\max \; r_{\text{IL}} + r_{\text{R}}\big(\phi_p(\quad)\big)$$

**Traffic Simulation Model**

{collision status,

distance to road boundary}

Realism Score

SMART-large

0.76

Behavior-GPT

SMART

max r (hand-
crafted features)

0.74

GUMP

0.72

IL pre-trained

Model Capacity (M)

0.70

1        5        10        >100

# Preference-Alignment with *Aligned* Features



*5000 preference rankings*

(A) ≻ (B) ≻ (C)

$\underset{\phi_R}{\operatorname{argmin}} \operatorname{diff}(\phi_R, \phi_H)$

$\phi_R \quad z \quad \phi_R^*$

$\phi_R^*(o^\pi)$
*Policy Feature Distribution*
⇕ matching
$\phi_R^*(o^+)$
*Preferred Feature Distribution*

$\max r_{IL} + r_R(\phi_R^*(\quad))$

**Traffic Simulation Model**

max r (aligned features)

SMART-large

**Behavior-GPT**

**SMART**

max r (hand-crafted features)

**GUMP**

**IL pre-trained**

**Realism Score**

0.76

0.74

0.72

0.70

**Model Capacity (M)**

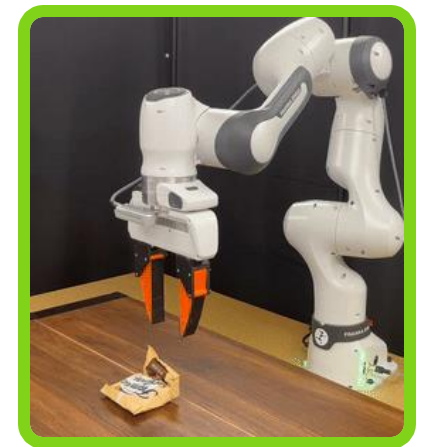1          5          10        >100

**Ours** is comparable to **100x larger** SOTA model

# Takeaways

**Model-based principles can help improve frontier robotics foundation models, but we need to "generalize" them to make them compatible with generative AI.**

① Robotics foundation model's training objective is only a proxy, we need post-training preference alignment

② We need robot representations to understand *what "matters" to us*

- Require **10x less** human budget to achieve high preference alignment in robotics manipulation

- Make a **1M** traffic model comparable to **100x larger** SOTA model after alignment with only a fraction of human data



**vs.**