# Lecture 9: Generative AI, Part 1

## Yilun Du

CS 2281R: Mathematical & Engineering Principles for
Training Foundation Models

# Overview

- Introduction to Multimodal Generative Models
- Compositional Generative Models

# Overview

- Introduction to Multimodal Generative Models

- **Compositional Generative Models**

# Generative Models of Language

Can you construct the steps for language based intelligence?

Creating language-based intelligence involves several structured steps that combine data collection, model development, evaluation, and deployment. Here's a high-level framework:

## 1. Data Collection and Preparation

- **Data Acquisition**: Gather large datasets of text relevant to your use case, like books, articles, transcriptions, or domain-specific texts.

- **Data Cleaning**: Remove any irrelevant content, duplicates, special characters, and format inconsistencies.

- **Data Labeling**: If the task requires it (like sentiment analysis or classification), label the data with the appropriate tags or annotations.

- **Tokenization**: Break down text into manageable units (words, subwords, or characters) depending on the model's requirements.
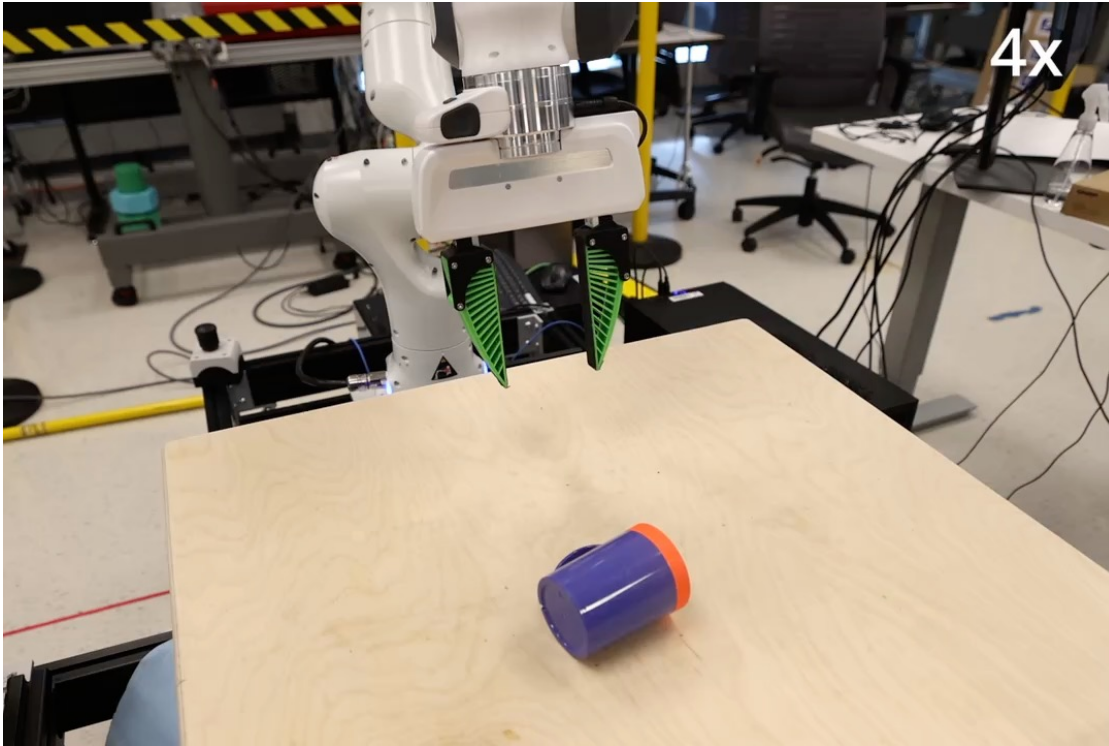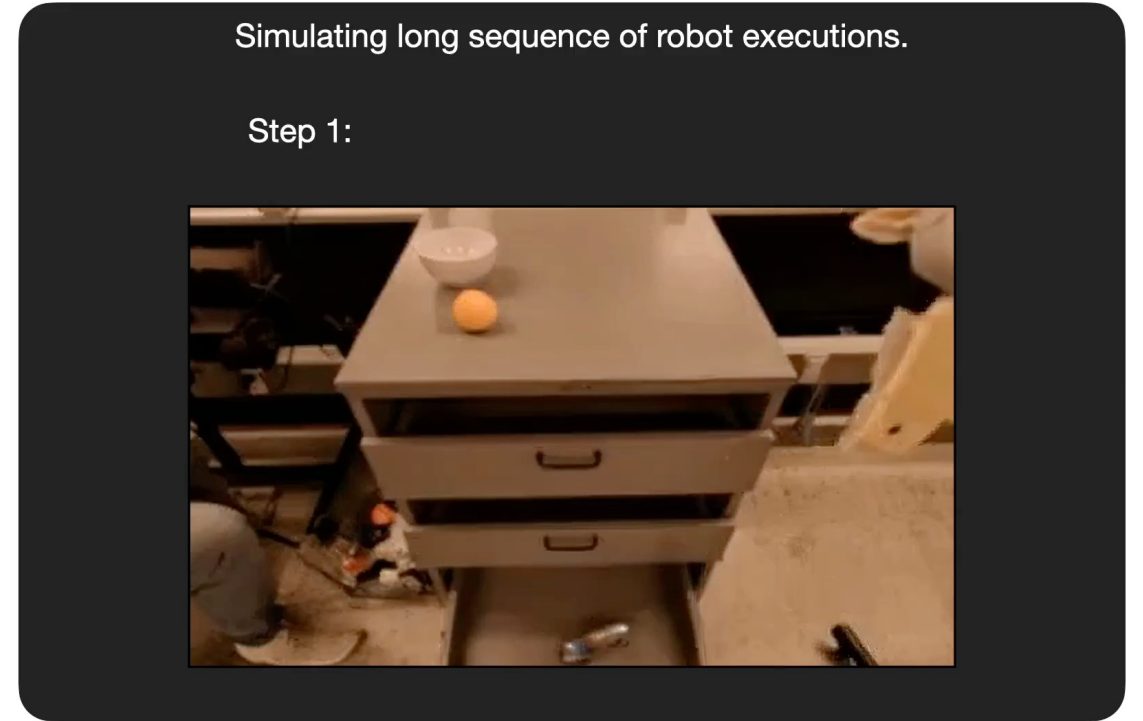
# Generative Models of Other Modalities



Images



Videos

# Generative Models of Other Modalities



Actions



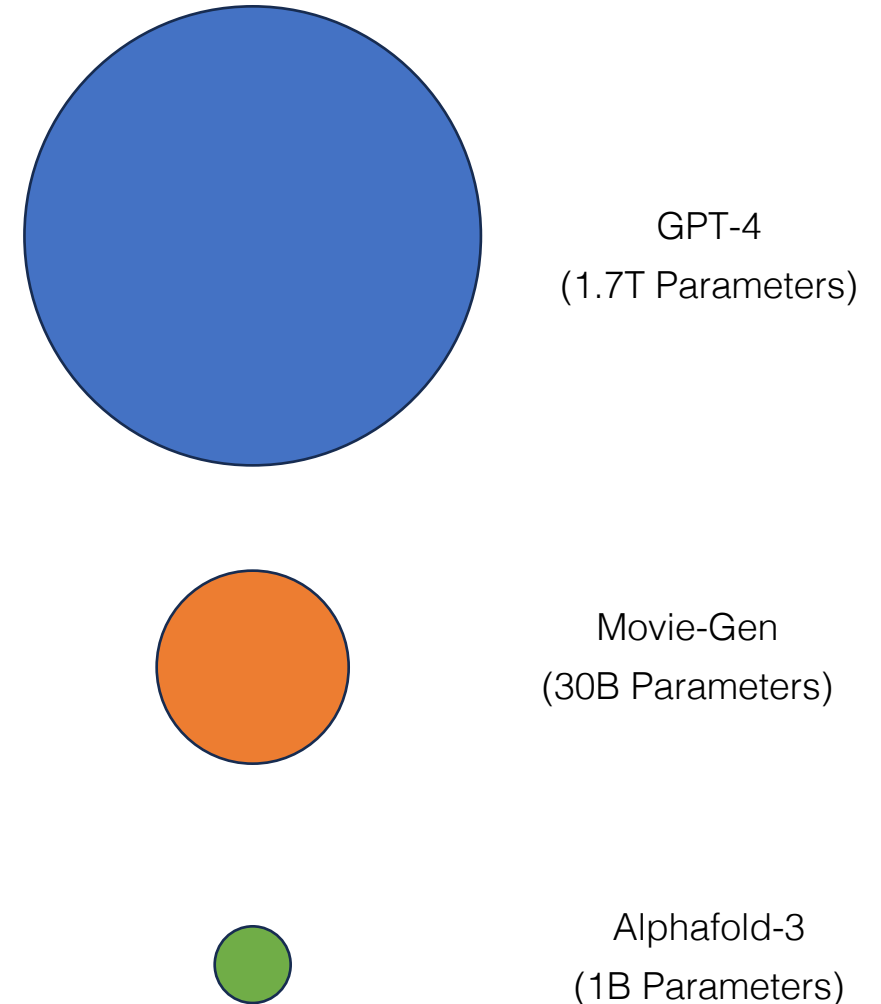Simulating long sequence of robot executions.

Step 1:

Simulation

# What are Some Challenges of Other Modalities?

- Individual variables in the distribution are not necessarily autoregressively dependent on each other….

- Distributions are much higher dimensional  than natural language, with much more uncertainty per pixel.

- We may not have data to cover the entire distribution we want to operate over.

# Relative Sizes of Models

- Are at a much smaller size than that of language models (but are trained on more data!).

- Scaling laws are much weaker than those seen in language.

- Models are very frail – text-to-image models often fail to follow even simple text prompts that deviate from those seen in training.

- Why?

  - Language may be uniquely information-rich and compositional…

GPT-4
(1.7T Parameters)

Movie-Gen
(30B Parameters)

Alphafold-3
(1B Parameters)

8

# A Tale of Computational Scaling



DCGAN (2015)



BigGAN (2018)



Stable Diffusion (2022)

# A Tale of Computational Scaling?



DCGAN (2015)



BigGAN (2018)

+ class conditioning
+ sample only high likelihood samples



Stable Diffusion (2022)

+ text conditioning
+ sample from distribution
$p(x)(p(x|c)/p(x))^\alpha$
+ very careful data filtering

# A Tale of Computational Scaling?



DCGAN (2015)

Model simple conditional distributions



BigGAN (2018)

+ class conditioning
+ sample only high likelihood samples



Stable Diffusion (2022)

+ text conditioning
+ sample from distribution $p(x)(p(x|c)/p(x))^\alpha$
+ very careful data filtering

# A Tale of Computational Scaling?



DCGAN (2015)

Generate samples from a
modified probability distribution



BigGAN (2018)

+ class conditioning
+ sample only high likelihood
samples



Stable Diffusion (2022)

+ text conditioning
+ sample from distribution
$p(x)(p(x|c)/p(x))^{\alpha}$
+ very careful data filtering

# A Tale of Computational Scaling?

An astronaut riding a horse



Stable Diffusion (2022)

+ text conditioning
+ sample from distribution
$p(x)(p(x|c)/p(x))^{\alpha}$
+ very careful data filtering

# A Tale of Computational Scaling?

An astronaut riding a horse



Stable Diffusion (2022)

+ text conditioning
+ sample from distribution
$p(x)(p(x|c)/p(x))^\alpha$
+ very careful data filtering

An astronaut riding a horse



Stable Diffusion (2022)

+ text conditioning
+ sample from distribution
$p(x)(p(x|c)/p(x))^\alpha$
+ very careful data filtering

# A Tale of Computational Scaling?

| An astronaut riding a horse | An astronaut riding a horse | A image |
|:---:|:---:|:---:|
|  |  |  |
| Stable Diffusion (2022) | Stable Diffusion (2022) | Stable Diffusion (2022) |

<div style="color:green">+ text conditioning</div>
<div style="color:green">+ sample from distribution</div>
$$p(x)(p(x|c)/p(x))^{\alpha}$$
<div style="color:green">+ very careful data filtering</div>

<div style="color:green">+ text conditioning</div>
<div style="color:red">+ sample from distribution</div>
$$p(x)(p(x|c)/p(x))^{\alpha}$$
<div style="color:green">+ very careful data filtering</div>

<div style="color:red">+ text conditioning</div>
<div style="color:red">+ sample from distribution</div>
$$p(x)(p(x|c)/p(x))^{\alpha}$$
<div style="color:green">+ very careful data filtering</div>

# A Tale of Computational Scaling?

An astronaut riding a horse



Stable Diffusion (2022)

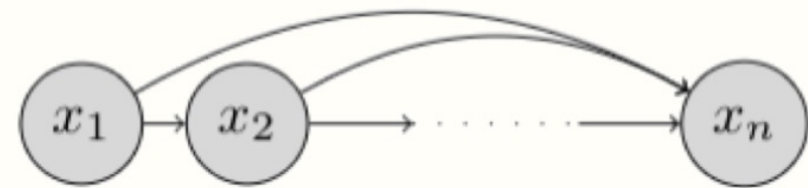An astronaut riding a horse



Stable Diffusion (2022)

A image



Stable Diffusion (2022)

Generative models cannot fit arbitrarily high dimensional distributions but rather ones that are simple (rich conditioning or low intrinsic dimensionality of data).

# Autoregressive Generative Models

- Language models are typically parameterized as autoregressive generative models
  - This reflects the natural casual order of language
- Can construct generative models over other distributions in an autoregressive manner
  - But this does not necessarily follow the "structure of the domain"
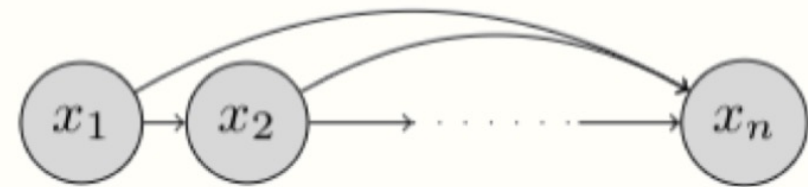  - In practice, we typically vector quantize continuous inputs into discrete tokens

$$p(\mathbf{x}) = \prod_{i=1}^{n} p(x_i | x_1, x_2, \ldots, x_{i-1}) = \prod_{i=1}^{n} p(x_i | \mathbf{x}_{<i})$$



17

# Limitations of Autoregressive Generative Models

- Learning an autoregressive factorization can be much harder than learning the base probability density itself:
  - For instance, consider a generative model of paths in a maze.

- The order of generation should follow the causal structure in the set of variables.
  - For example, any-order language models perform poorly
  - Generate pixels of an image in a hierarchical manner

$$p(\mathbf{x}) = \prod_{i=1}^{n} p(x_i | x_1, x_2, \dots, x_{i-1}) = \prod_{i=1}^{n} p(x_i | \mathbf{x}_{<i})$$

$x_1 \rightarrow x_2 \rightarrow \cdots \cdots \rightarrow x_n$

# Other Possible Generative Models

- There is a zoo of other generative models such as:
  - Energy-based Models
  - Variational Autoencoders
  - GANs
  - Flow Models
  - Diffusion models
  - Many of these generative model classes can be interconverted between each other

# Energy Based Models

- The oldest class of generative models that inspired the development of many of the generative models we know today
  - Noise Contrastive Estimation -> GANs, Variational Partition Function Minimization -> VAEs, Ancestral Importance Sample + Score Matching -> Diffusion
- Represent the probability distribution $p_\theta(x)$ as an unnormalized distribution parameterized with an energy function $E_\theta(x)$ where $p_\theta(x) \propto e^{-E_\theta(x)}$
  - Allows us to represent any probability distribution with a neural network that outputs a single scalar output
- The simplest training objective for EBMs is approximate maximum likelihood estimation

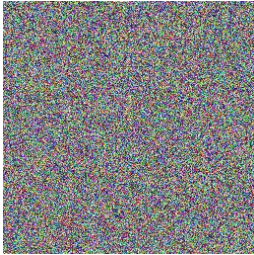# Energy Based Models

- The likelihood of datapoint x is given under an EBM is give by $p_\theta(x) = e^{-E\theta(x)} / \int e^{-E\theta(x)} dx$ , where denominator is known as the partition function and is usually intractable to compute.

- The gradient of maximum likelihood training for a point $x$ is given by:

$$\nabla_\theta \log p_\theta(x) = -\nabla_\theta E_\theta(\text{x}) - \nabla_\theta \log \int e^{-E\theta(x)} dx$$

$$= -\nabla_\theta E_\theta(\text{x}) + \frac{\int \nabla_\theta E_\theta(x) e^{-E\theta(x)} dx}{\int e^{-E\theta(x)} dx}$$

$$= -\nabla_\theta E_\theta(\text{x}) + \mathbb{E}_{x \sim p_\theta(x)}[\nabla_\theta E_\theta(x)]$$

- The last expression can estimate through Monte Carlo approximation by sampling from the model distribution!

# Energy Based Models

- The likelihood of datapoint x is given under an EBM is give by $p_\theta(x) = e^{-E_\theta(x)} / \int e^{-E_\theta(x)} dx$ , where denominator is known as the partition function and is usually intractable to compute.
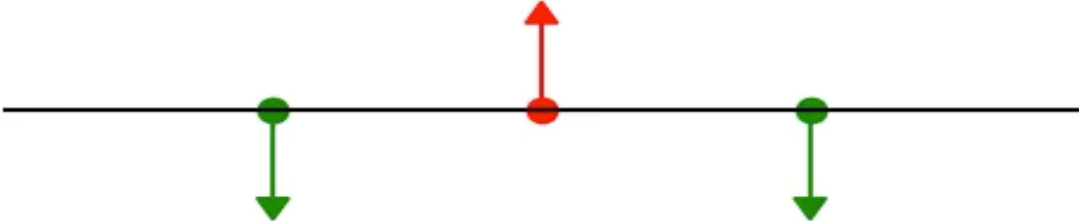
- This corresponds to maximum likelihood loss:

$$\nabla_\theta \mathcal{L}_{ML} = \underbrace{\mathbb{E}_{x \sim p(x)}}[\nabla_\theta E(x)] - \underbrace{\mathbb{E}_{x \sim p_\theta(x)}}[\nabla_\theta E(x)]$$

Drawn from Data     Drawn from Learned Distribution

- Contrastively decrease the energy of real data while increase the energy of samples from the model (inspired the development of contrastive learning).
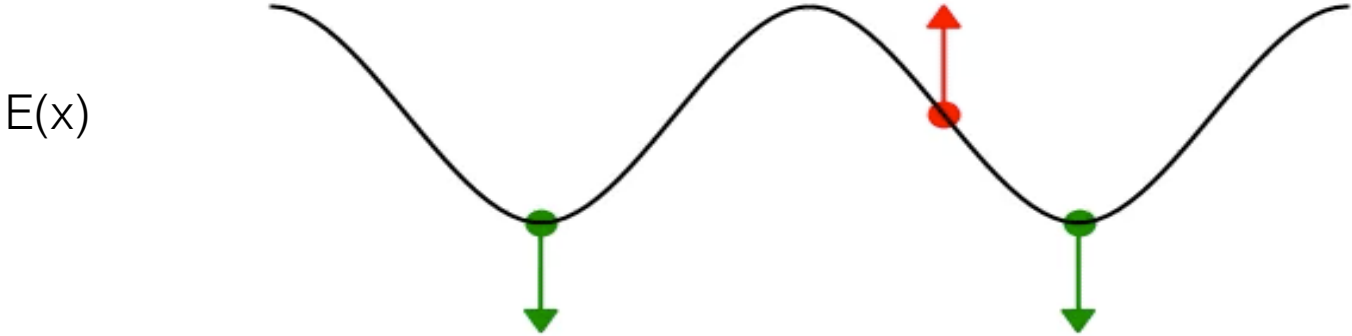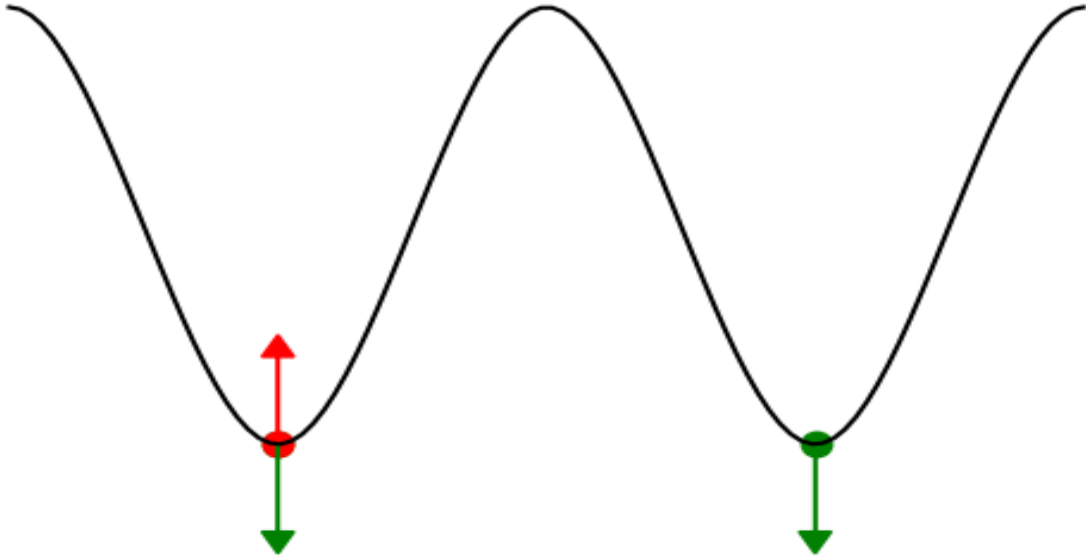
# Learning Energy Functions



E(x)

# Learning Energy Functions



E(x)

# Learning Energy Functions



$E(x)$

# Limitations of Energy Based Models

- EBMs make training very challenging and slow because it requires you to explicitly draw MCMC sample from the model probability distribution in order to train the model
  - This inspired many future generative models which learned explicit networks for sampling
- However, this comes with benefits. – the parameterization $p_\theta(x) \propto e^{-E_\theta(x)}$ makes no assumptions on the nature of the probability distribution modeled, allowing the model to learn to flexibly capture any distribution
- In contrast, all other generative models make assumptions on the structure of the probability distribution they are modeling, which can be inaccurate dependent on the distribution.

# From EBMs to Variational Autoencoders

- Training EBMs is challenging because it involves sampling from a high dimensional distribution $p_\theta(x)$ for maximum likelihood training

- We can make it easier by factorizing the distribution with a latent distribution

$$p_\theta(x) = \int p_\theta(x|z)p(z)dz$$

- Learning $p_\theta(x|z)$ can be much easier than learning $p_\theta(x)$ -- for instance $p_\theta(x|z)$ can be Gaussian even when $p_\theta(x)$ is not

- However, maximum likelihood training of $p_\theta(x)$ still requires us to exhaustively sample all value of $p(z)$ which is intractable

# From EBMs to Variational Autoencoders

- Use variational inference to learn an amortized sampler $q_\emptyset(z|x)$ for $p(z)$ given an input $x$

$$p_\theta(x) = \int p_\theta(x|z)p(z)dz = \int p_\theta(x|z)p(z)\frac{q_\emptyset(z|x)}{q_\emptyset(z|x)}dz = \text{E}_{q_\emptyset(x|z)}\left[p(z)\frac{p_\theta(x|z)}{q_\emptyset(z|x)}\right]$$

- Using the Jensen's inequality, we can write the log-likelihood as:

$$\log p_\theta(x) = \log\left(\text{E}_{q_\emptyset(z|x)}\left[p(z)\frac{p_\theta(x|z)}{q_\emptyset(x|z)}\right]\right) \geq \text{E}_{q_\emptyset(z|x)}[p_\theta(x|z)] - KL(q_\emptyset(z|x)||p(z))$$

- The above objective is exactly the training objective used to train the VAE!

# Variational Autoencoders

- Maximize log-likelihood in the form

$$\log p_\theta(x) = \log \left( \mathrm{E}_{q_\emptyset(z|x)} \left[ p(z) \frac{p_\theta(x|z)}{q_\emptyset(x|z)} \right] \right) \geq \mathrm{E}_{q_\emptyset(z|x)}[p_\theta(x|z)] - KL(q_\emptyset(z|x)||p(z))$$

- In a VAE, we represent $p_\theta(x|z)$ and $q_\emptyset(z|x)$ as Gaussian distributions

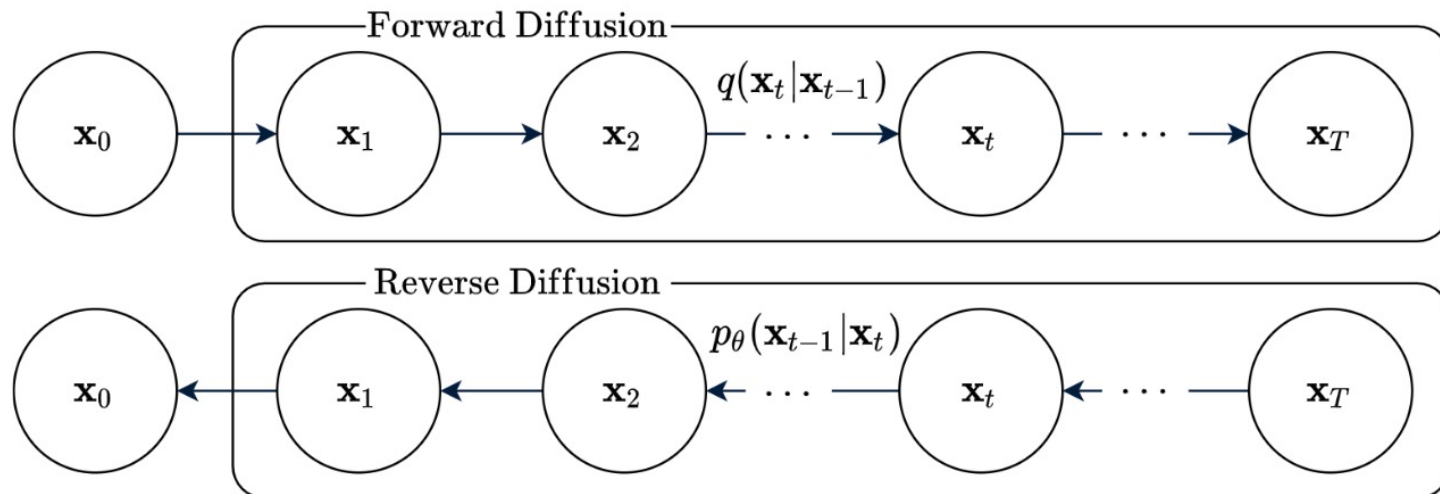# From Variational Autoencoders / EBMs to Diffusion Models

- In practice, VAEs often generate blurry samples as both amortized sampler and generator have limited capacity
- Reduce the capacity requirements for each component of the variational procedure by constructing an annealed sequence of intermediate latent variables (inspired from anneal importance sampling from EBMs)

$$p\left(\mathbf{x}^{(0\cdots T)}\right) = p\left(\mathbf{x}^{(T)}\right) \prod_{t=1}^{T} p\left(\mathbf{x}^{(t-1)}|\mathbf{x}^{(t)}\right).$$

$$\log \, p\left(\mathbf{x}^{(0\cdots T)}\right) = \log \left[ \begin{array}{c} \int d\mathbf{x}^{(1\cdots T)} q\left(\mathbf{x}^{(1\cdots T)}|\mathbf{x}^{(0)}\right) \cdot \\ p\left(\mathbf{x}^{(T)}\right) \prod_{t=1}^{T} \frac{p\left(\mathbf{x}^{(t-1)}|\mathbf{x}^{(t)}\right)}{q\left(\mathbf{x}^{(t)}|\mathbf{x}^{(t-1)}\right)} \end{array} \right]$$

# Diffusion Models

- In practice, VAEs often generate blurry samples as both amortized sampler and generator have limited capacity
- Reduce the capacity requirements for each component of the variational procedure by constructing an annealed sequence of intermediate latent variables (developed originally to draw samples from the partition function of EBMs)
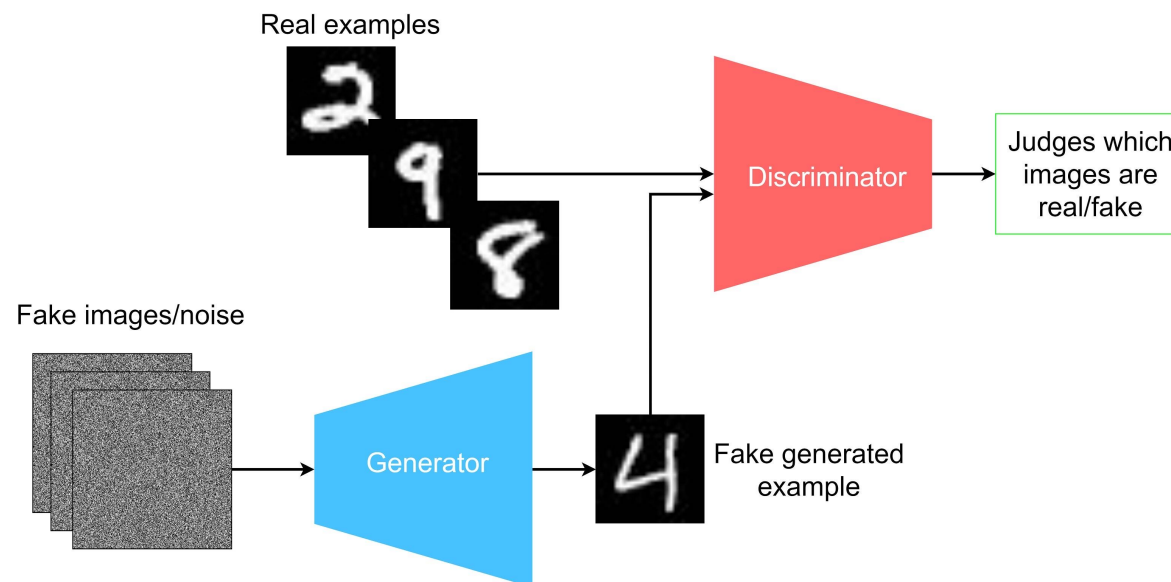
$$\mathbb{E}\left[-\log p_\theta(\mathbf{x}_0)\right] \le \mathbb{E}_q\left[-\log \frac{p_\theta(\mathbf{x}_{0:T})}{q(\mathbf{x}_{1:T}|\mathbf{x}_0)}\right] = \mathbb{E}_q\left[-\log p(\mathbf{x}_T) - \sum_{t\ge 1}\log \frac{p_\theta(\mathbf{x}_{t-1}|\mathbf{x}_t)}{q(\mathbf{x}_t|\mathbf{x}_{t-1})}\right] =: L$$

$$\mathbb{E}_q\left[\underbrace{D_{\mathrm{KL}}(q(\mathbf{x}_T|\mathbf{x}_0) \,\|\, p(\mathbf{x}_T))}_{L_T} + \sum_{t>1}\underbrace{D_{\mathrm{KL}}(q(\mathbf{x}_{t-1}|\mathbf{x}_t,\mathbf{x}_0) \,\|\, p_\theta(\mathbf{x}_{t-1}|\mathbf{x}_t))}_{L_{t-1}} \underbrace{-\log p_\theta(\mathbf{x}_0|\mathbf{x}_1)}_{L_0}\right]$$

# Diffusion Models

- In practice, VAEs often generate blurry samples as both amortized sampler and generator have limited capacity

- Reduce the capacity requirements for each component of the variational procedure by constructing an annealed sequence of intermediate latent variables (developed originally to draw samples from the partition function of EBMs)

# From EBMs to GANs

- To learn an EBM $p_\theta(x) \propto e^{-E_\theta(x)}$ to fit a probability distribution $p_D(x)$, one clever way is through classification (noise contrastive estimation)

- Given samples from a data distribution $p_D(x)$ and a noise distribution $p_{noise}(x)$, we can implicitly recover the energy function $E_\theta(x)$ by training a classifier $\mathcal{H}(x)$ classifying if a data point is either from the data or noise distribution. The energy function is the difference of the logits $\mathcal{H}(x)$ and $p_{noise}(x)$.

- This procedure allows us to replace the difficult problem of drawing samples from model distribution with generating samples from the noise distribution. However, the variance of this estimator still depends on how close $p_{noise}(x)$ is to sampling from the $p_\theta(x)$

# GANs

- Instead of explicitly constructing a distribution $p_{noise}(x)$, learn the neural network g(z) that approximates this distribution!

- Continue using the classifier $\mathcal{H}(x)$ to distinguish between real and fake samples

- The noise distribution is the generator and the classifier is the discriminator, leading to GANs!

Real examples

Discriminator

Judges which images are real/fake

Fake images/noise

Generator

Fake generated example

34

# Overview

- Introduction to Multimodal Generative Models
- **Compositional Generative Models**

# Compositional Generative Models

- In practice, in many multimodal domains, data is scarce and often covers a sparse subset of the distribution we would like generative models to operate over.

- We can use the idea of compositional generative modeling to induce models to generate samples that are outside of the distribution of the data seen by the model

$$p(x, y, z) = p(x, y) \, p(y, z)$$

- Above factorization only requires data to be seen from the pairwise marginals instead of the full joint distribution of samples

# Why Does Generative AI for Language Work So Well?

Real World Distribution

Natural Language

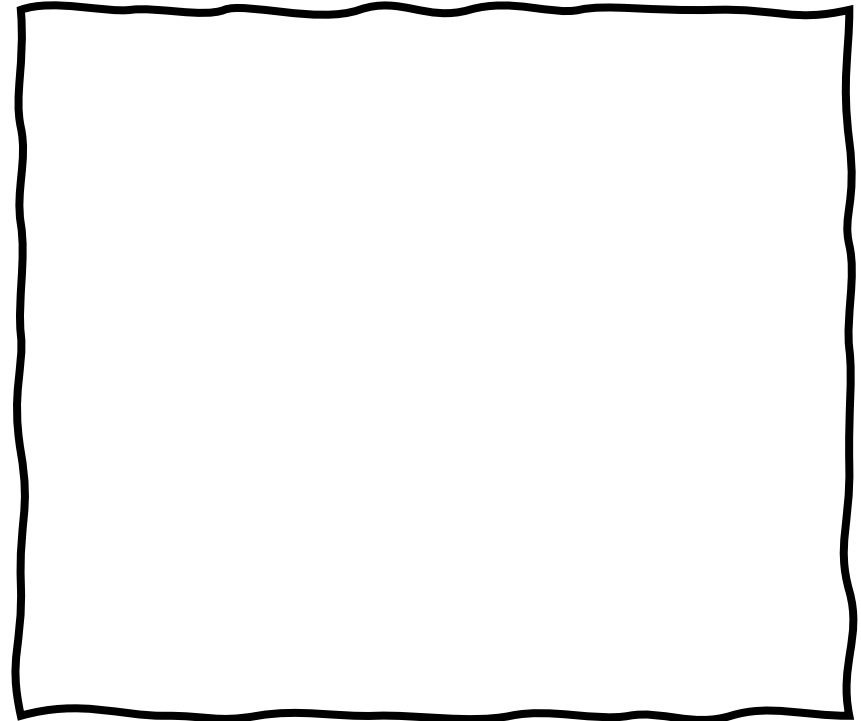# Why Does Generative AI for Language Work So Well?

Real World Distribution

Training Distribution

Natural Language

# Why Does Generative AI on Other Settings Work Poorly?
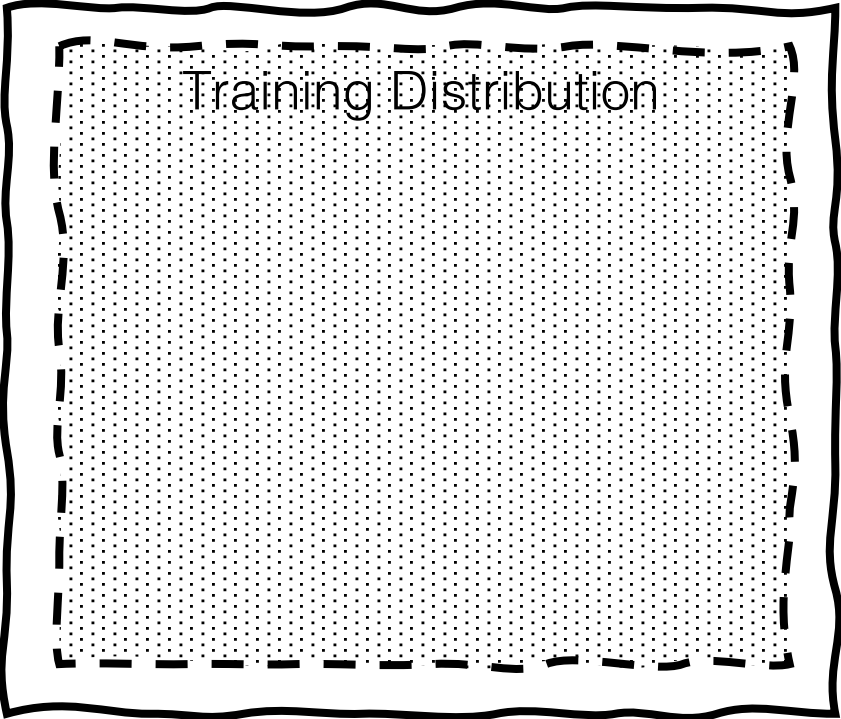
Real World Distribution

Training Distribution

Natural Language

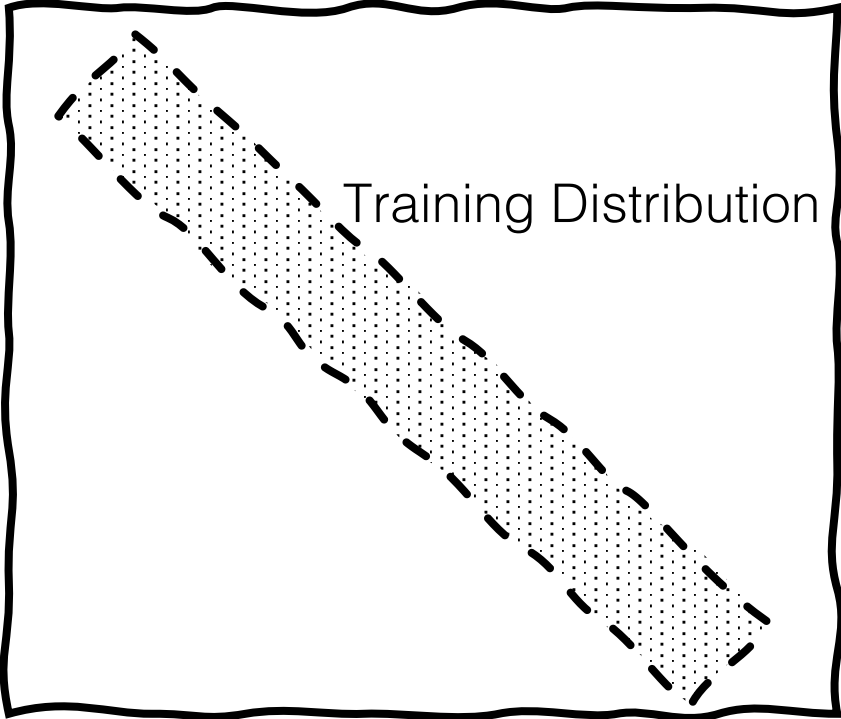Real World Distribution

Embodied Data

# Why Does Generative AI on Other Settings Work Poorly?
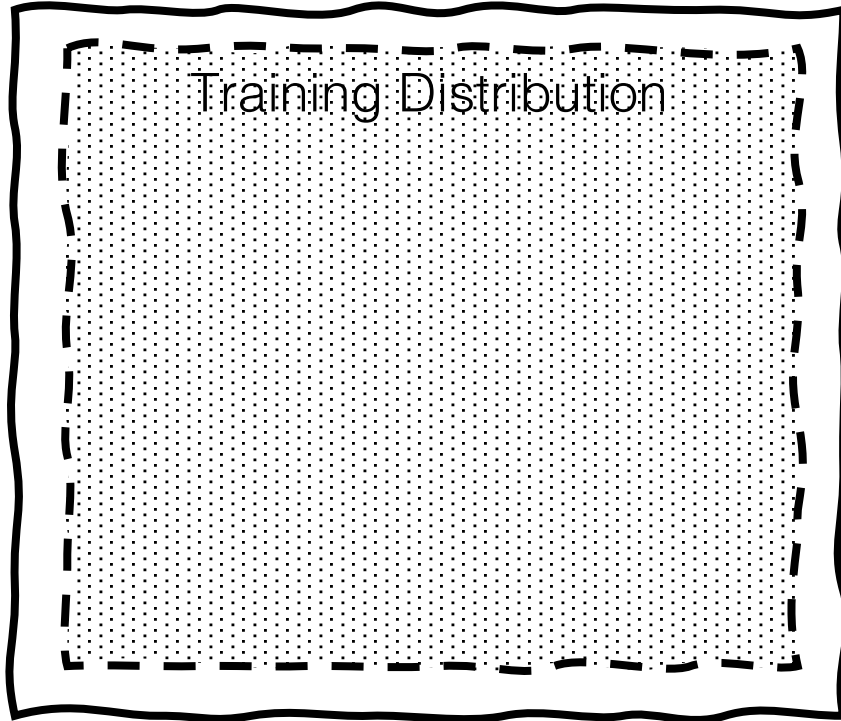
Real World Distribution

Training Distribution

Natural Language

Real World Distribution

Training Distribution

Embodied Data

# Why Does Generative AI on Other Settings Work Poorly?
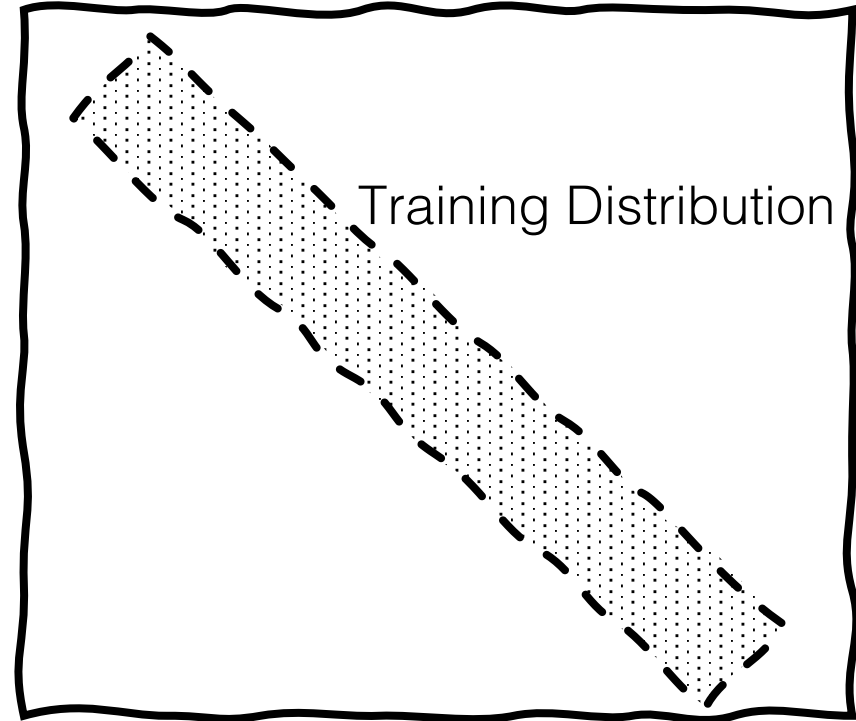
Real World Distribution

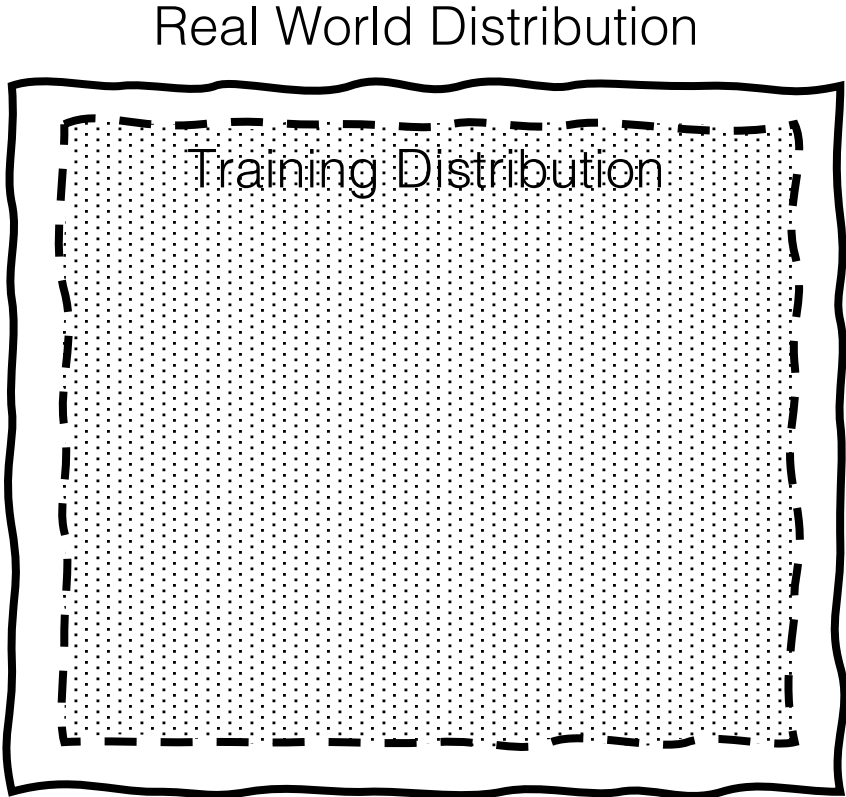Training Distribution

Natural Language

Real World Distribution
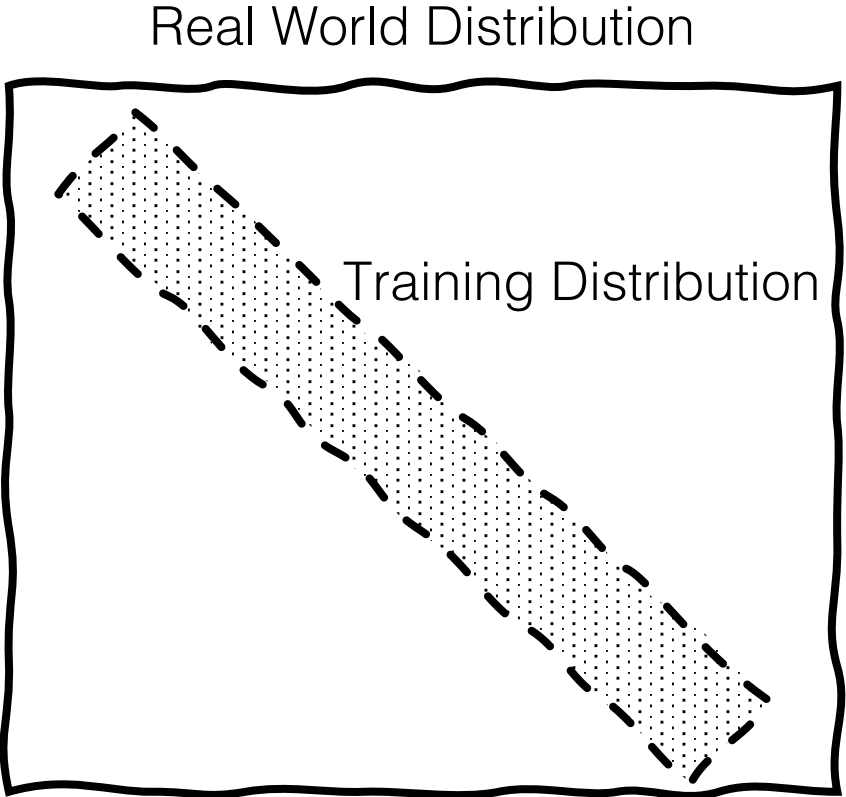
Training Distribution

Embodied Data

Compared to pixels in images (and many other continuous distributions), natural language is much simpler (structured for effective communication) and naturally compositional (infinite use of finite means).

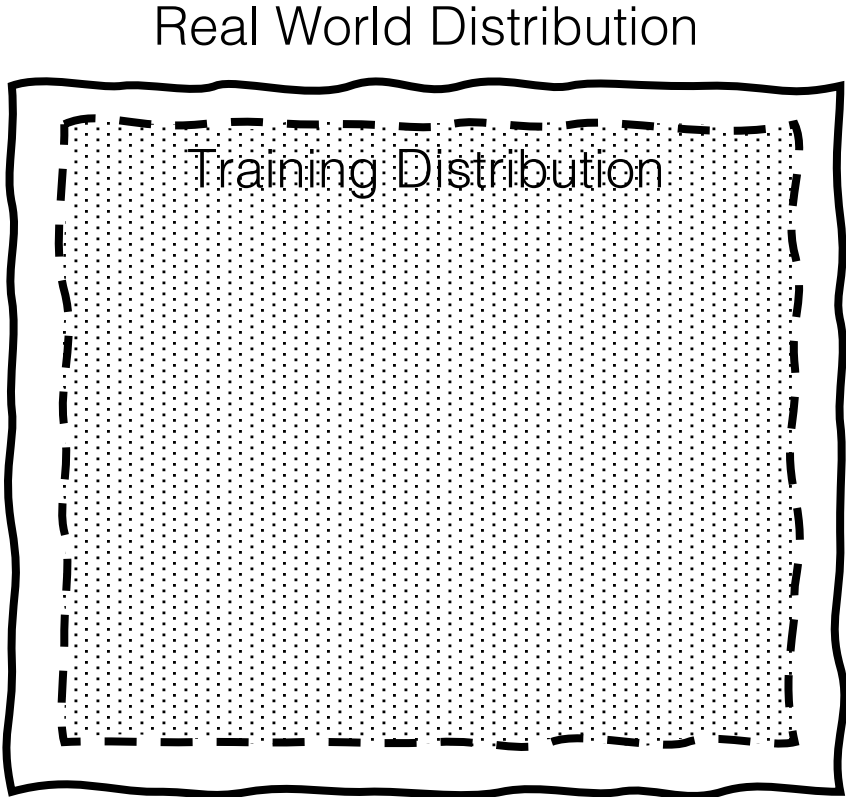# Composition of Learned Factors Yields Strong Generalization

Real World Distribution

Training Distribution

Natural Language

Real World Distribution

Training Distribution

Embodied Data

Energy Based Models (EBMs) provide a probabilistic manner to represent the real distribution as a composition of factors!

# Composition of Learned Factors Yields Strong Generalization

Real World Distribution



Training Distribution

Natural Language

Real World Distribution



$$p(x) \propto \prod_{i=1}^{N} p_i(x_i)$$

Training Distribution

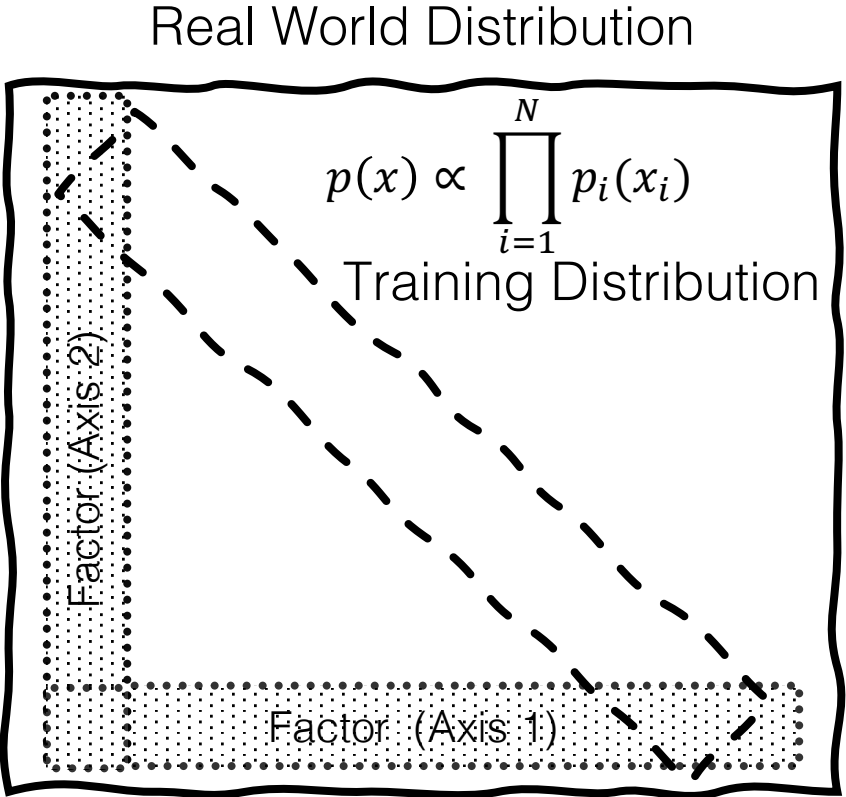Factor (Axis 2)

Factor (Axis 1)

Embodied Data

Energy Based Models (EBMs) provide a probabilistic manner to represent the real distribution as a composition of factors!

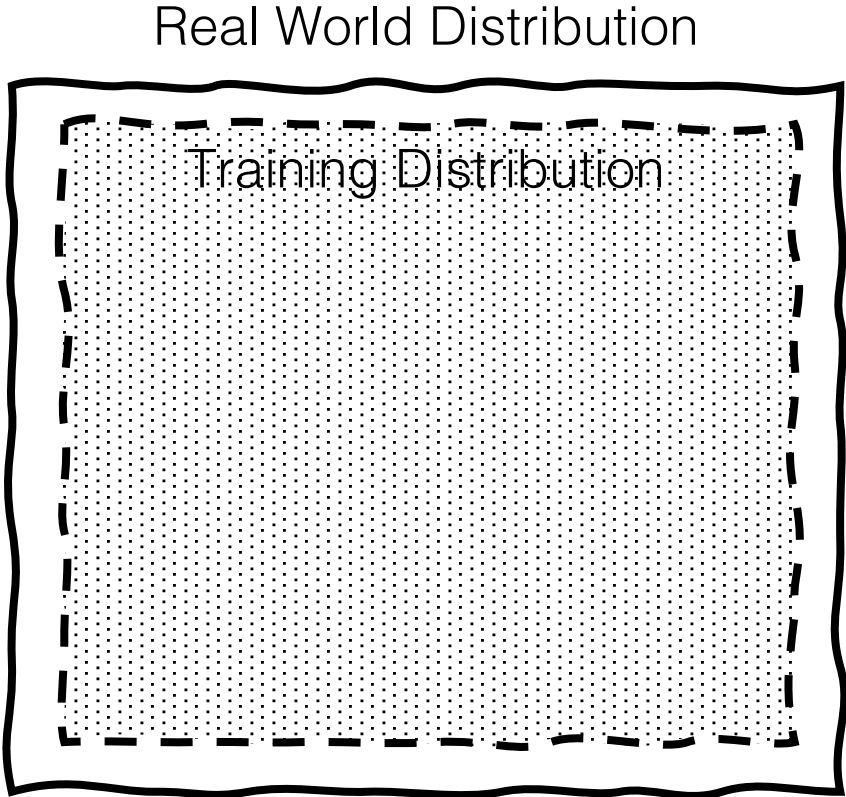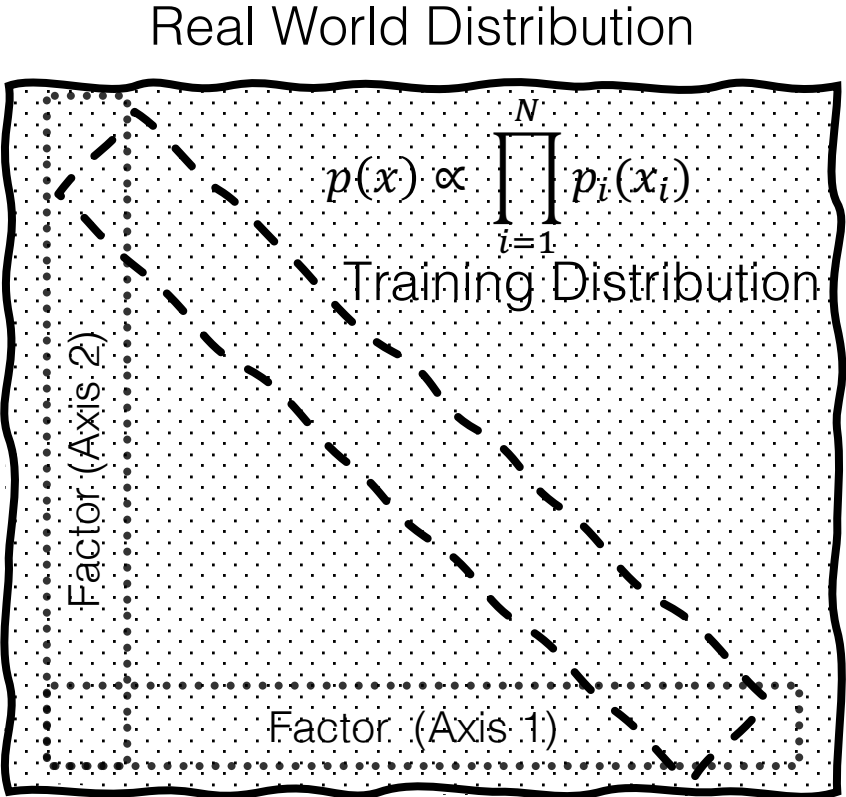# Composition of Learned Factors Yields Strong Generalization

Real World Distribution

Training Distribution

Natural Language

Real World Distribution

$$p(x) \propto \prod_{i=1}^{N} p_i(x_i)$$

Training Distribution

Factor (Axis 2)

Factor (Axis 1)
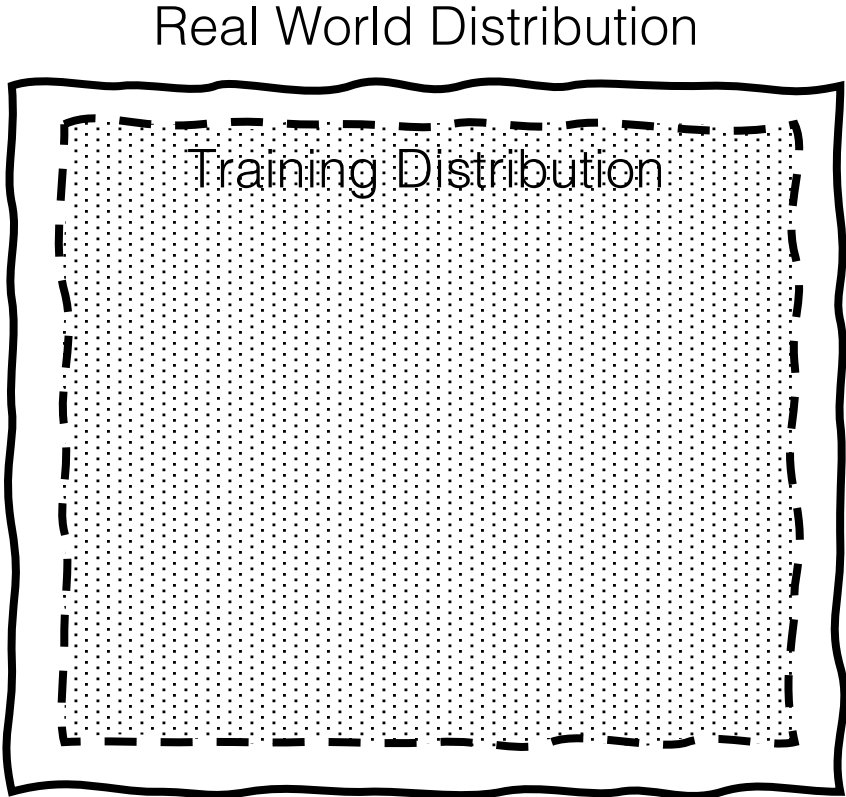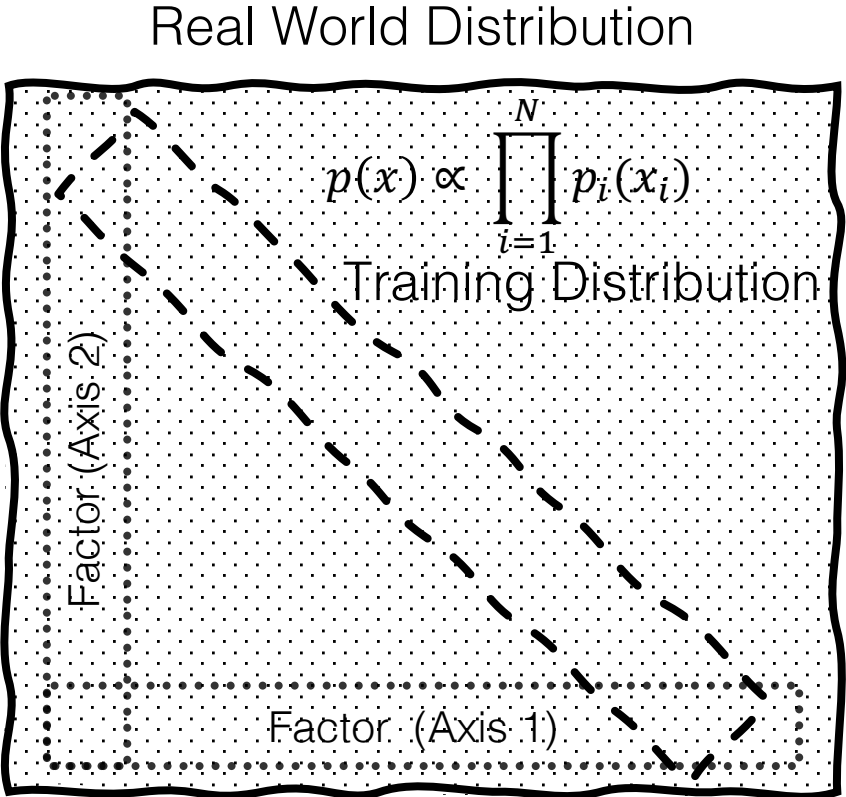
Embodied Data

Factors are combined to represent the entire distribution (even parts with no data).

# Composition of Learned Factors Yields Strong Generalization



Real World Distribution

Training Distribution

Natural Language

Real World Distribution

$$p(x) \propto \prod_{i=1}^{N} p_i(x_i)$$
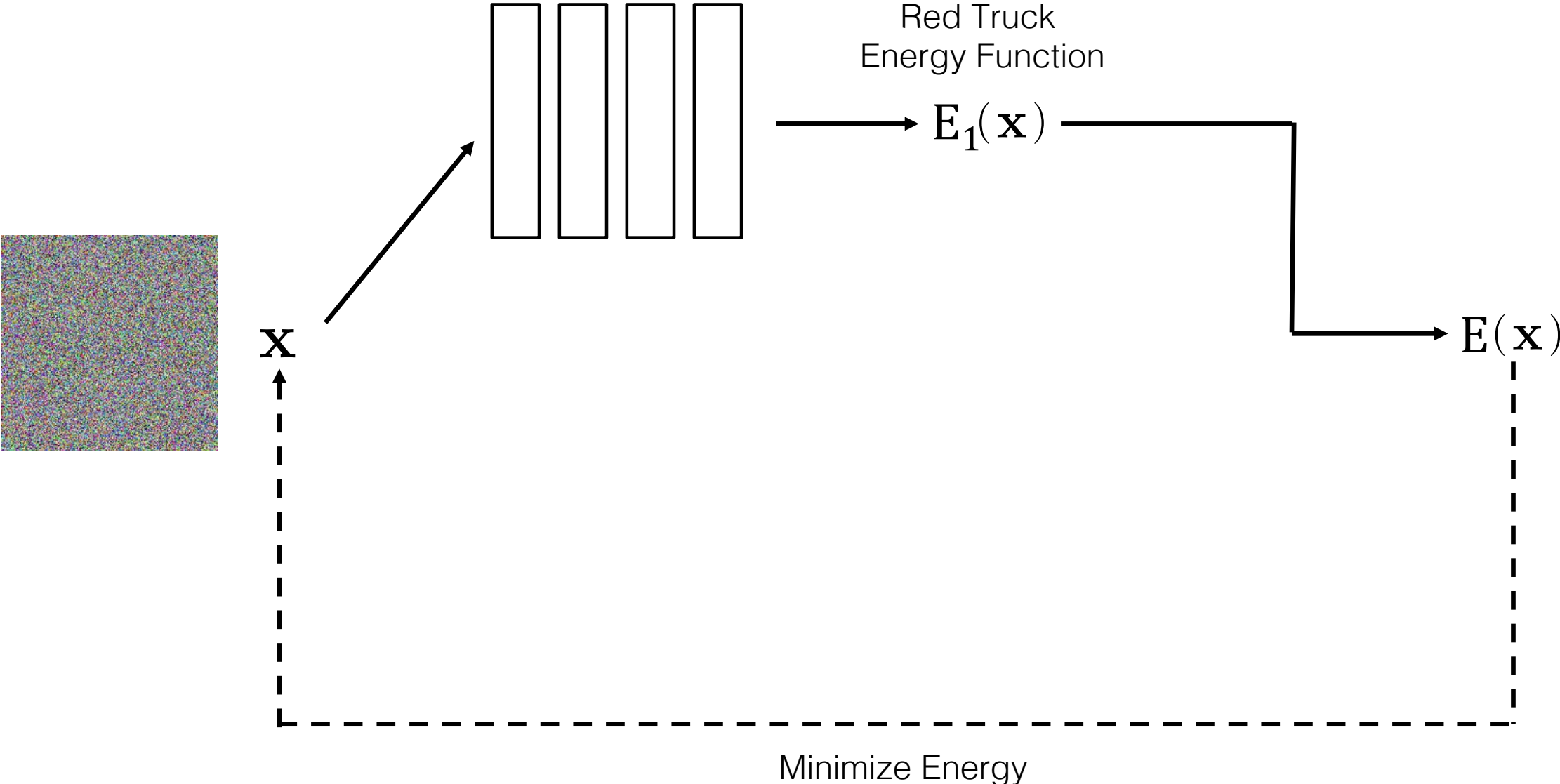
Training Distribution

Factor (Axis 2)

Factor (Axis 1)

Embodied Data

Factors make independence assumptions about data that are biased.

# Composing Different Energy Functions at Prediction Time



Red Truck
Energy Function

$$E_1(\mathbf{x})$$

$$\mathbf{x}$$

$$E(\mathbf{x})$$

Minimize Energy

# Composing Different Energy Functions at Prediction Time



Red Truck
Energy Function

$$E_1(\mathbf{x})$$

$$\mathbf{x}$$

Desert
Energy Function

$$E_n(\mathbf{x})$$

$$E(\mathbf{x})$$

Minimize Energy

# Combining Probability Distributions with Energy Functions

Products:



$p_1(x)$ $\times$ $p_2(x)$ $=$ $p_1(x)\, p_2(x)$

[1] Du et al. Reduce, Reuse, Recycle: Compositional Generation with Energy-Based Diffusion Models and MCMC. ICML 2023

# Combining Probability Distributions with Energy Functions

Products:



$p_1(x)$ × $p_2(x)$ = $p_1(x)\, p_2(x)$

Mixtures:



$p_1(x)$ + $p_2(x)$ = $p_1(x) + p_2(x)$

[1] Du et al. Reduce, Reuse, Recycle: Compositional Generation with Energy-Based Diffusion Models and MCMC. ICML 2023

# Combining Probability Distributions with Energy Functions



Products:

$p_1(x)$ × $p_2(x)$ = $p_1(x)\, p_2(x)$

Mixtures:

$p_1(x)$ + $p_2(x)$ = $p_1(x) + p_2(x)$

Inversion:

$p_1(x)$ / $p_2(x)$ = $p_1(x)\, /\, p_2(x)^a$

[1] **Du** et al. Reduce, Reuse, Recycle: Compositional Generation with Energy-Based Diffusion Models and MCMC. ICML 2023

# Applications of Composing Generative Models



Single Model

Compositional Model

"A couch right next to the windows" AND "A table in front of the couch" AND "A vase of flowers on top of the table"

"A blue bird on a tree" AND "A red car behind the tree" AND "A green forest in the background"

"A green tree swaying in the wind" AND "A red brick house located behind a tree" AND "A lawn in front of the house"

"A pink sky" AND "A blue mountain in the horizon" AND "Cherry Blossoms in front of the mountain"

[1] Liu*, Li*, **Du\*** et al. Composable Visual Generation with Diffusion Models. ECCV 2022

Original Generation

Adapted Generation (Digital Art)

Adapted Generation (Outdoor Video)

Adapted Generation (Storybook Illustration)

[2] Yang*, **Du\*** et al. Probabilistic Adaptation of Text-to-Video Models. ICLR 2024

Movie still of epic space battle

Starship Enterprise firing phasers

Giant mocha robot holding a glowing sword

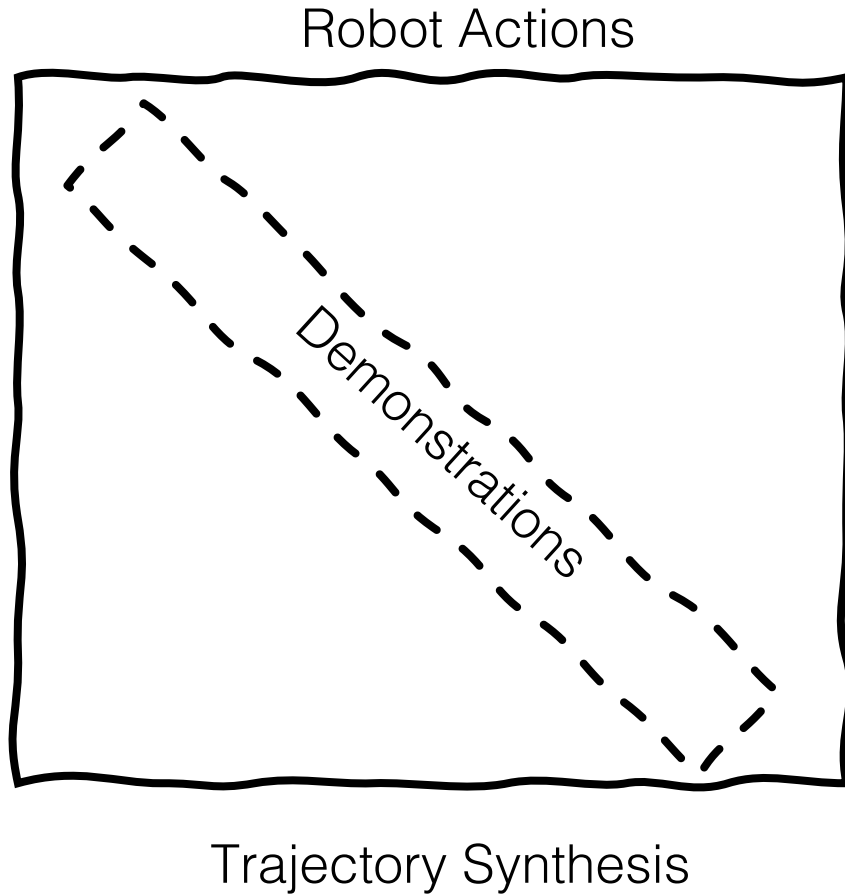Glowing phaser beam

Sun with lens flare

Portion of Mars. Debris from atmosphere

[3] **Du** et al. Reduce, Reuse, Recycle: Compositional Generation with Energy-Based Diffusion Models and MCMC. ICML 2023

$geom_{Box}$   $geom_A$   $geom_B$   $geom_C$

$in(A, Box)$   $close\text{-}to(A, B)$   $cfree(A, B)$   $cfree(A, C)$ ......

$pose_{Box}$   $pose_A$   $pose_B$   $pose_C$

$pose_{A0}$

$grasp_A$   $valid\text{-}traj$   $traj_A$

The arm trajectory $traj_A$ connects A's initial pose $pose_{A0}$ and the target pose $pose_A$ given $grasp_A$.

Box   A   D   C   E   B

(a) Visualization of the environment while placing object A.

(b) Visualization of the constraint graphs associated with the object placement. There are three decision variables.
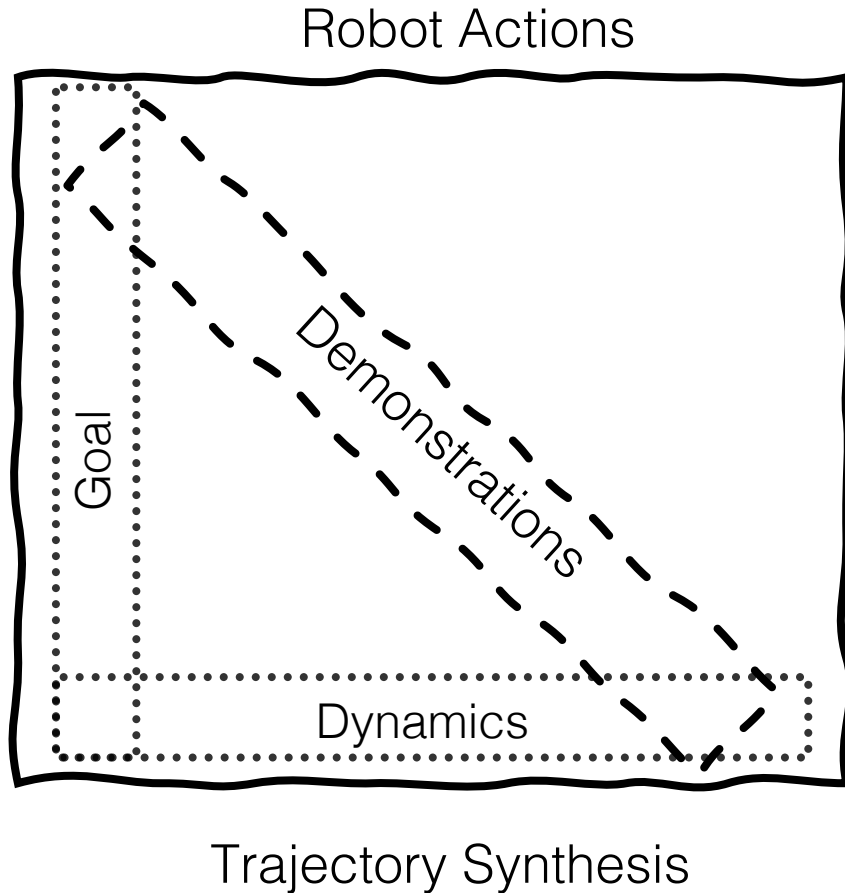
[4] Yang, Mao, **Du** et al. Compositional Diffusion-Based Continuous Constraint Solvers. CoRL 2023

# Generalizing Beyond Demonstrations through Compoistion

Robot Actions

Demonstrations

Trajectory Synthesis

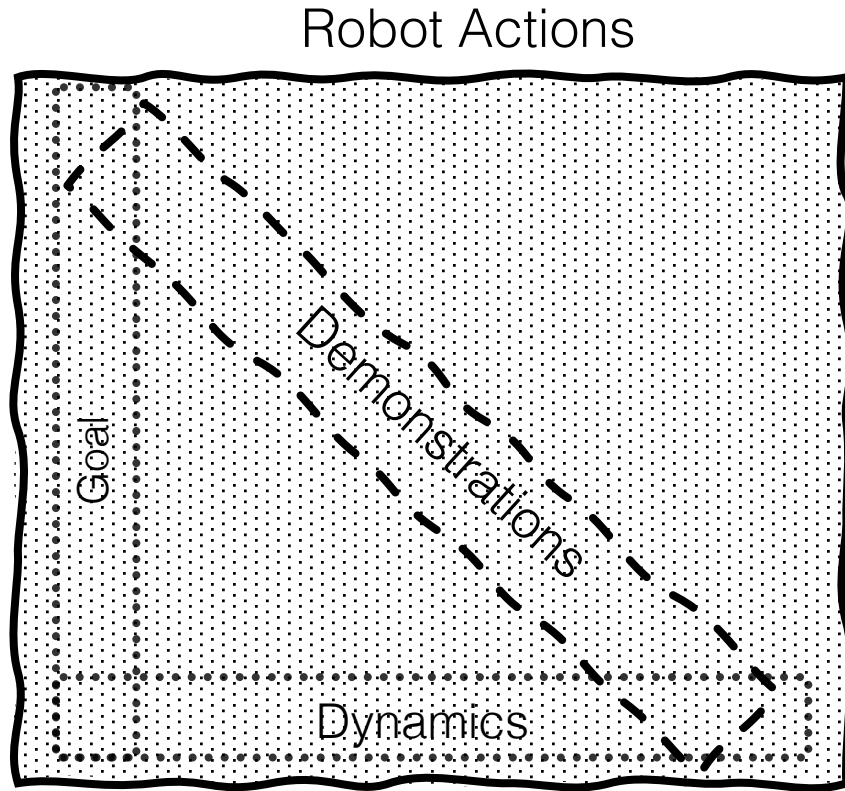We want to construct robotic agents that can generalize beyond the demonstrations they have seen!

# Planning through Compositional Generation

Robot Actions

Goal

Demonstrations

Dynamics

Trajectory Synthesis

Decompose demonstrations into a model capturing the dynamics of the environment + a model to capture the goal of a task.
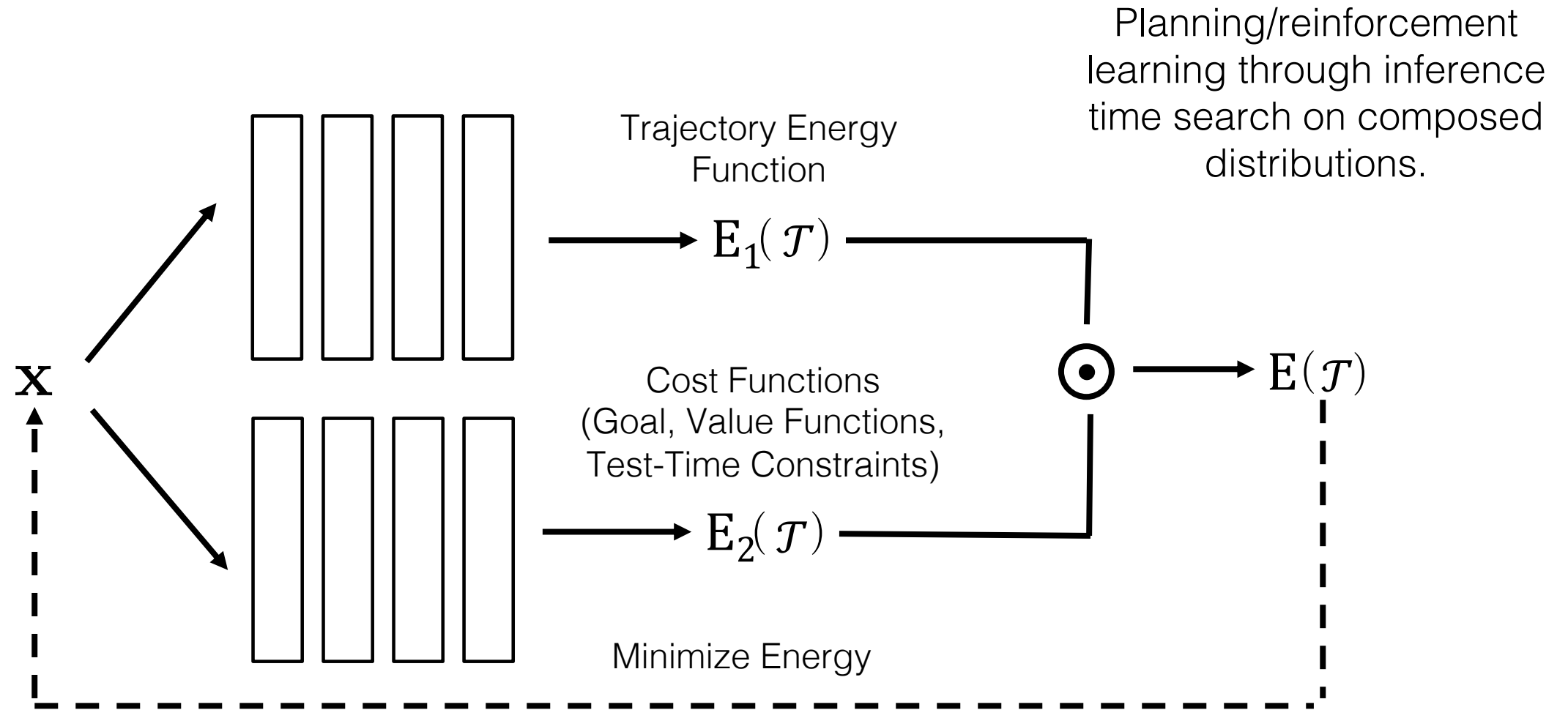
$$p_{\mathrm{plan}}(\tau) \propto p_{\mathrm{dynamics}}(\tau) p_{\mathrm{task}}(\tau, \mathrm{goal})$$

# Planning through Compositional Generation

Robot Actions

Goal

Demonstrations

Dynamics

Trajectory Synthesis

Enables generalization to new combinations of states + goals through probabilistic planning!

$$p_{\text{plan}}(\tau) \propto p_{\text{dynamics}}(\tau) p_{\text{task}}(\tau, \text{goal})$$

54

# Planning with Energy Minimization



Trajectory Energy Function

$$\mathbf{E}_1(\mathcal{T})$$

Cost Functions
(Goal, Value Functions,
Test-Time Constraints)

$$\mathbf{E}_2(\mathcal{T})$$

Minimize Energy

Planning/reinforcement learning through inference time search on composed distributions.

$$\mathbf{E}(\mathcal{T})$$

[1] **Du** et al. Model Based Planning with Energy Based Models. CoRL 2019.
[2] Janner*, **Du*** et al. Planning with Diffusion for Flexible Behavior Synthesis. ICML 2022
[3] Ajay*, **Du***, Gupta* et al. Is Conditional Generative Modeling all You Need For Decision Making? ICLR 2023
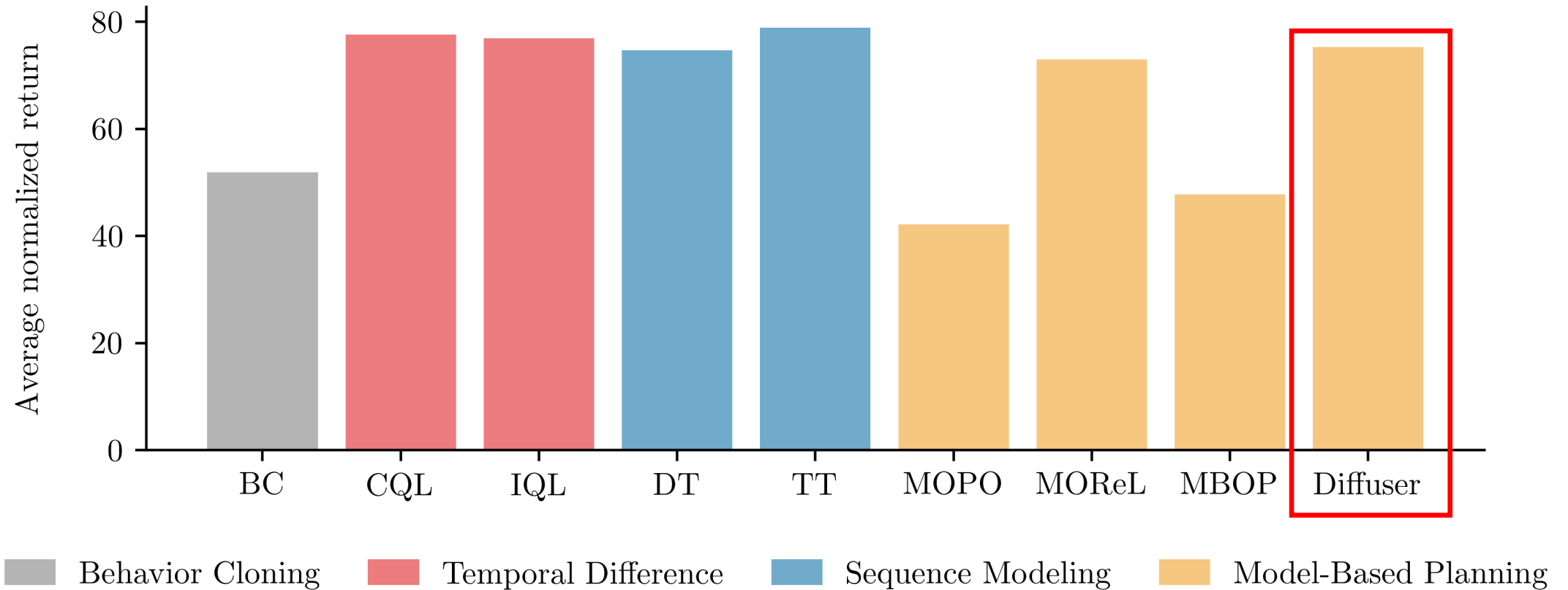
# Reinforcement Learning through Value Composition

$E_1(\mathcal{T})$

$\mathbf{s}_0 \quad \mathbf{s}_1 \quad \mathbf{s}_2 \quad \mathbf{s}_3 \quad \mathbf{s}_4 \quad \mathbf{s}_5$

$\mathbf{a}_0 \quad \mathbf{a}_1 \quad \mathbf{a}_2 \quad \mathbf{a}_3 \quad \mathbf{a}_4 \quad \mathbf{a}_5$

$E_2(\mathcal{T})$

$r(\mathbf{s}_0) + \gamma r(\mathbf{s}_1) + \gamma^2 r(\mathbf{s}_2) + \gamma^3 r(\mathbf{s}_3) + \gamma^4 r(\mathbf{s}_4) + \gamma^5 V(\mathbf{s}_5)$

$\oplus \longrightarrow E_1(\mathcal{T}) + E_2(\mathcal{T})$

# Reinforcement Learning through Value Composition

$E_1(\mathcal{T})$

$$\mathbf{s}_0 \quad \mathbf{s}_1 \quad \mathbf{s}_2 \quad \mathbf{s}_3 \quad \mathbf{s}_4 \quad \mathbf{s}_5$$

$$\mathbf{a}_0 \quad \mathbf{a}_1 \quad \mathbf{a}_2 \quad \mathbf{a}_3 \quad \mathbf{a}_4 \quad \mathbf{a}_5$$

$E_2(\mathcal{T})$

$$r(\mathbf{s}_0) + \gamma r(\mathbf{s}_1) + \gamma^2 r(\mathbf{s}_2) + \gamma^3 r(\mathbf{s}_3) + \gamma^4 r(\mathbf{s}_4) + \gamma^5 V(\mathbf{s}_5)$$

$\bigoplus \longrightarrow E_1(\mathcal{T}) + E_2(\mathcal{T})$

Solve many different tasks with one model!

# Reinforcement Learning through Value Composition

# Goal Planning through Optimization



$E_1(\mathcal{T})$

$$s_0 \quad s_1 \quad s_2 \quad s_3 \quad s_4 \quad s_5$$

$$a_0 \quad a_1 \quad a_2 \quad a_3 \quad a_4 \quad a_5$$

$E_2(\mathcal{T})$

$$s_0 \qquad\qquad\qquad\qquad g$$

$\oplus \longrightarrow E_1(\mathcal{T}) + E_2(\mathcal{T})$

$$s_0 \quad s_1 \quad s_2 \quad s_3 \quad s_4 \quad g$$

$$a_0 \quad a_1 \quad a_2 \quad a_3 \quad a_4 \quad a_5$$

Construct a zero-shot goal-seeking policy!

# Goal Planning through Optimization

# Goal Planning through Optimization

Baselines require per task retraining
while our trained model can be applied
across tasks!



Single Task Planning

Multi-Task Planning

# Planning From Partial Visual Observations on Real Robots



Can learn complex trajectory planning given only visual observations given very few (50) demos.

[1] Chi, Feng, **Du** et al. Diffusion Policy: Visualmotor Policy Learning via Action Diffusion. RSS 2023.

# Solving Long Horizon Tasks by Composing Foundation Models



Make A Cup of Tea

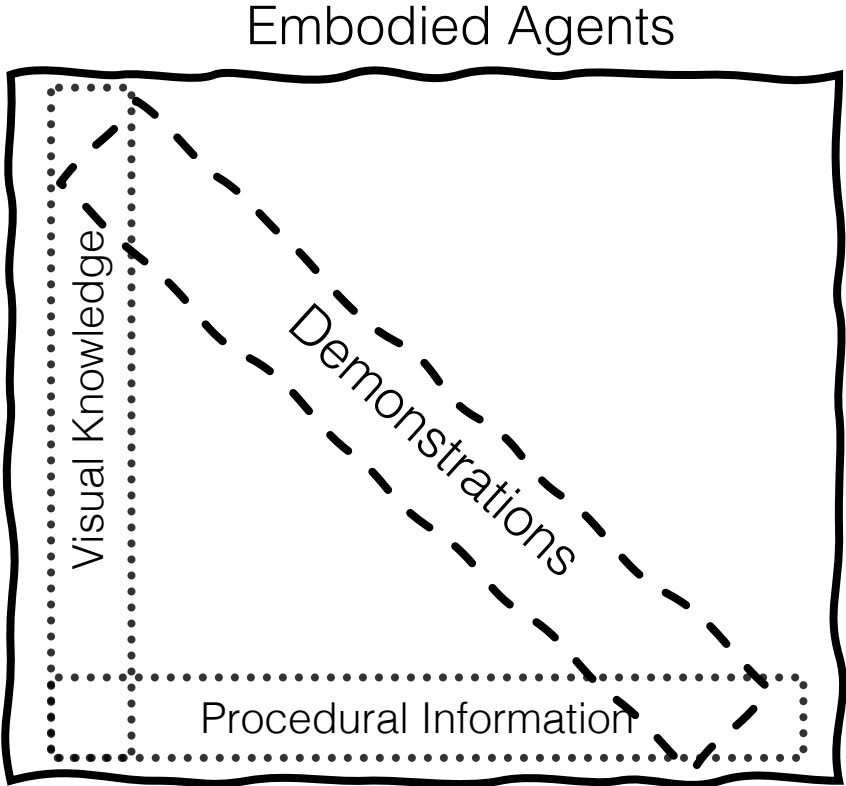Solving a long horizon decision-making task requires generalization across many sources of knowledge.

Embodied Agents

Demonstrations

Decision Making

# Composing Foundation Models for Embodied Agents



Make A Cup of Tea

Compose foundation models representing each axis of information!



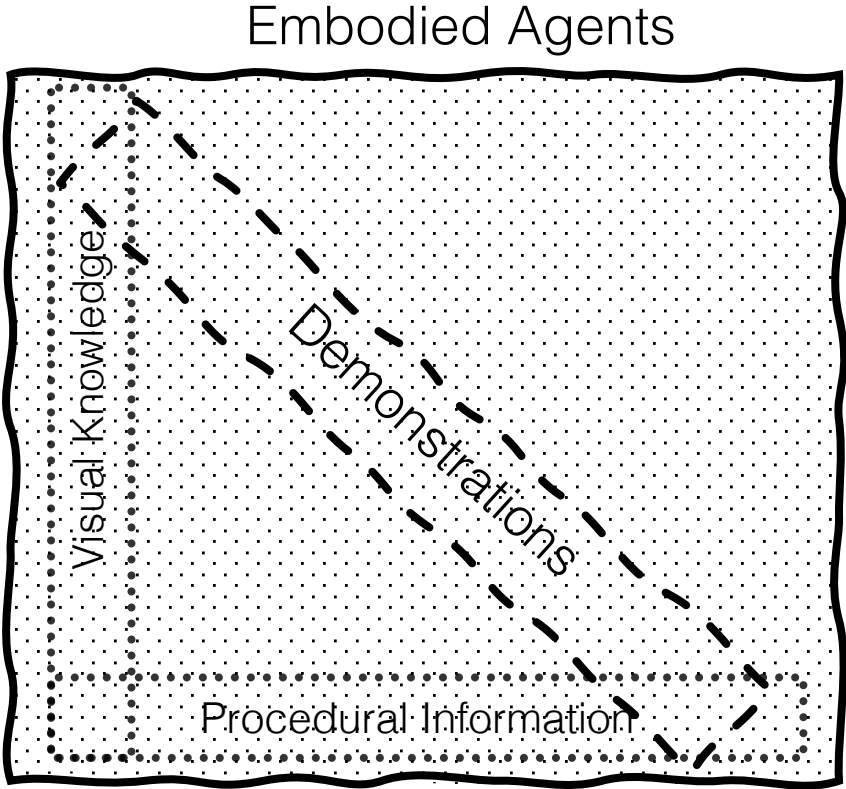Embodied Agents

Visual Knowledge

Demonstrations

Procedural Information

Decision Making

$$p_{\text{tea}}(\tau_{\text{text}}, \tau_{\text{image}}, \tau_{\text{action}}) \propto$$
$$p_{\text{LLM}}(\tau_{\text{text}}) \, p_{\text{video}}(\tau_{\text{text}}, \tau_{\text{image}}) \, p_{\text{action}}(\tau_{\text{image}}, \tau_{\text{action}})$$

# Composing Foundation Models for Embodied Agents
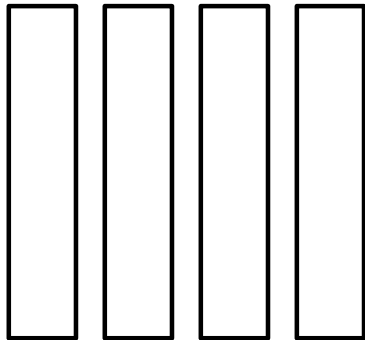


Make A Cup of Tea

Compose foundation models representing each axis of information!

Embodied Agents

Decision Making

$$p_{\text{tea}}(\tau_{\text{text}}, \tau_{\text{image}}, \tau_{\text{action}}) \propto$$
$$p_{\text{LLM}}(\tau_{\text{text}}) \, p_{\text{video}}(\tau_{\text{text}}, \tau_{\text{image}}) \, p_{\text{action}}(\tau_{\text{image}}, \tau_{\text{action}})$$
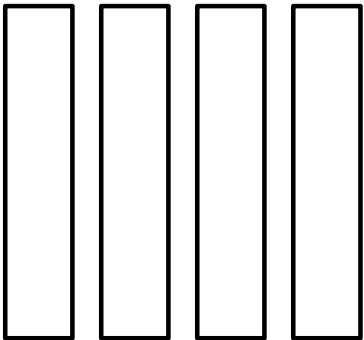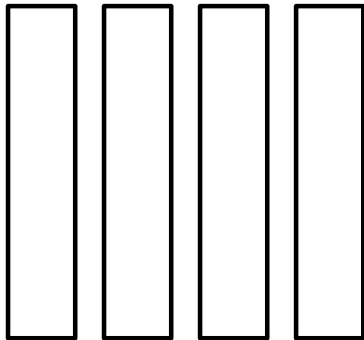
# Hierarchical Planning with Foundation Models

1) Look for a tea kettle
2) Heat water
3) Find teabag
4) …



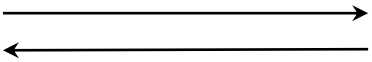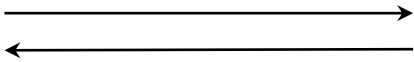Task Information

Motion Information

Kinematics Information

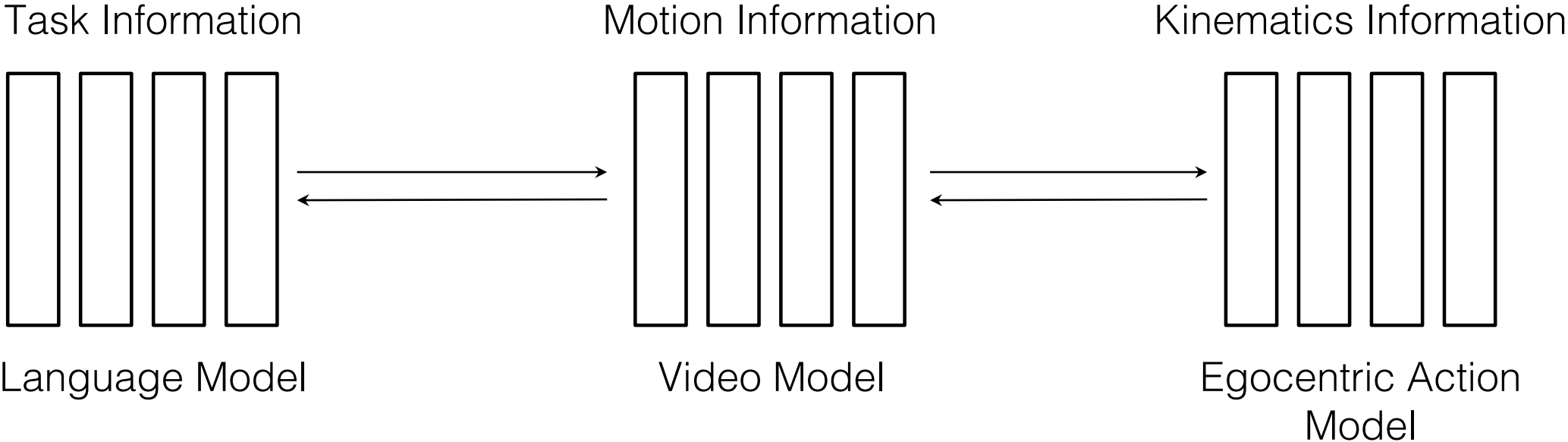Language Model

Video Model

Egocentric Action Model

We want to construct a plan to make a cup of tea that is semantically, geometrically, and physically executable on a robot!

[1] Ajay*, Han*, **Du*** et al. Compositional Foundation Models for Hierarchical Planning. NeurIPS 2023

# Hierarchical Planning with Foundation Models

$$p_{\text{tea}}(\tau_{\text{text}}, \tau_{\text{image}}, \tau_{\text{action}}) \propto$$
$$p_{\text{LLM}}(\tau_{\text{text}})\, p_{\text{video}}(\tau_{\text{text}}, \tau_{\text{image}})\, p_{\text{action}}(\tau_{\text{image}}, \tau_{\text{action}})$$
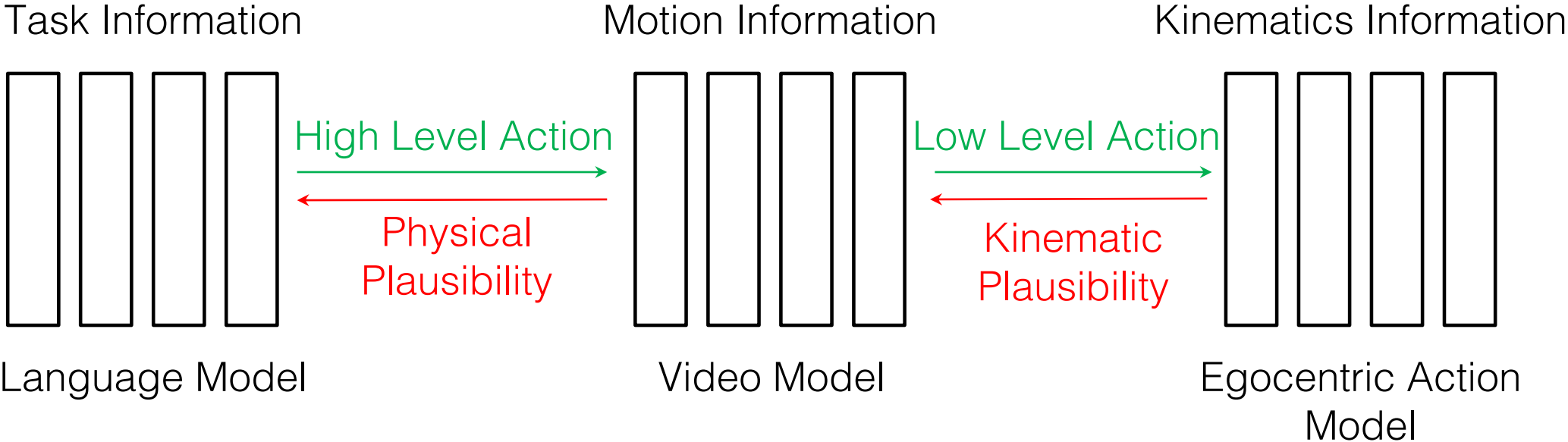
Task Information       Motion Information       Kinematics Information



Language Model       Video Model       Egocentric Action Model

We want to construct a plan to make a cup of tea that is
semantically, geometrically, and physically executable on a robot!

[1] Ajay*, Han*, **Du\*** et al. Compositional Foundation Models for Hierarchical Planning. NeurIPS 2023

# Hierarchical Planning with Foundation Models

$$p_{\text{tea}}\big(\tau_{\text{text}}, \tau_{\text{image}}, \tau_{\text{action}}\big) \propto$$
$$p_{\text{LLM}}(\tau_{\text{text}})\, p_{\text{video}}\big(\tau_{\text{text}}, \tau_{\text{image}}\big)\, p_{\text{action}}\big(\tau_{\text{image}}, \tau_{\text{action}}\big)$$

Task Information                     Motion Information                     Kinematics Information

High Level Action                    Low Level Action

Physical Plausibility                Kinematic Plausibility

Language Model                       Video Model                          Egocentric Action Model

We want to construct a plan to make a cup of tea that is
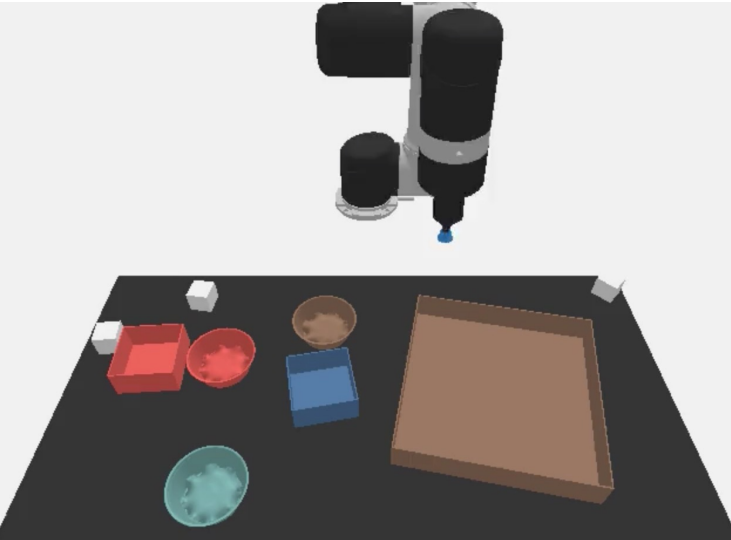semantically, geometrically, and physically executable on a robot!

[1] Ajay*, Han*, **Du*** et al. Compositional Foundation Models for Hierarchical Planning. NeurIPS 2023

# Hierarchical Decision Making with Multimodal Models

## Goal

Stack red block on a cyan block and place a brown block to the right of stack.

## Start Image



## Generated Language Plan

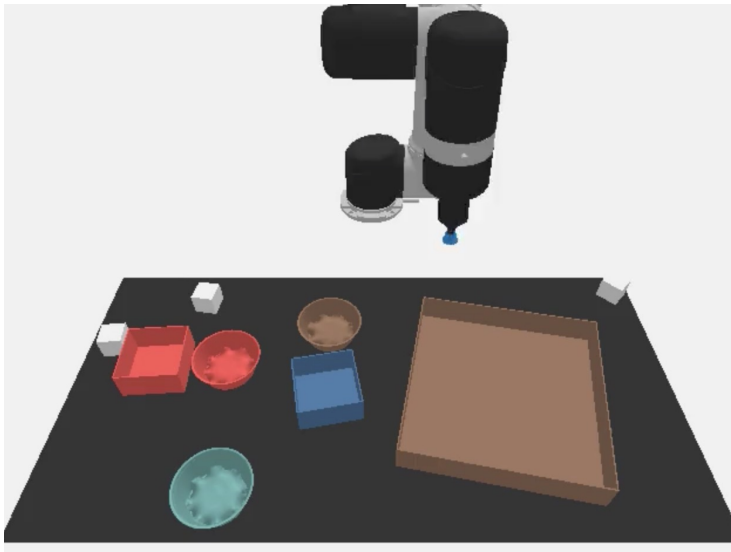Place cyan block in brown box.

## Generated Video Plan

# Hierarchical Decision Making with Multimodal Models

## Goal

Stack red block on a cyan block and place a brown block to the right of stack.
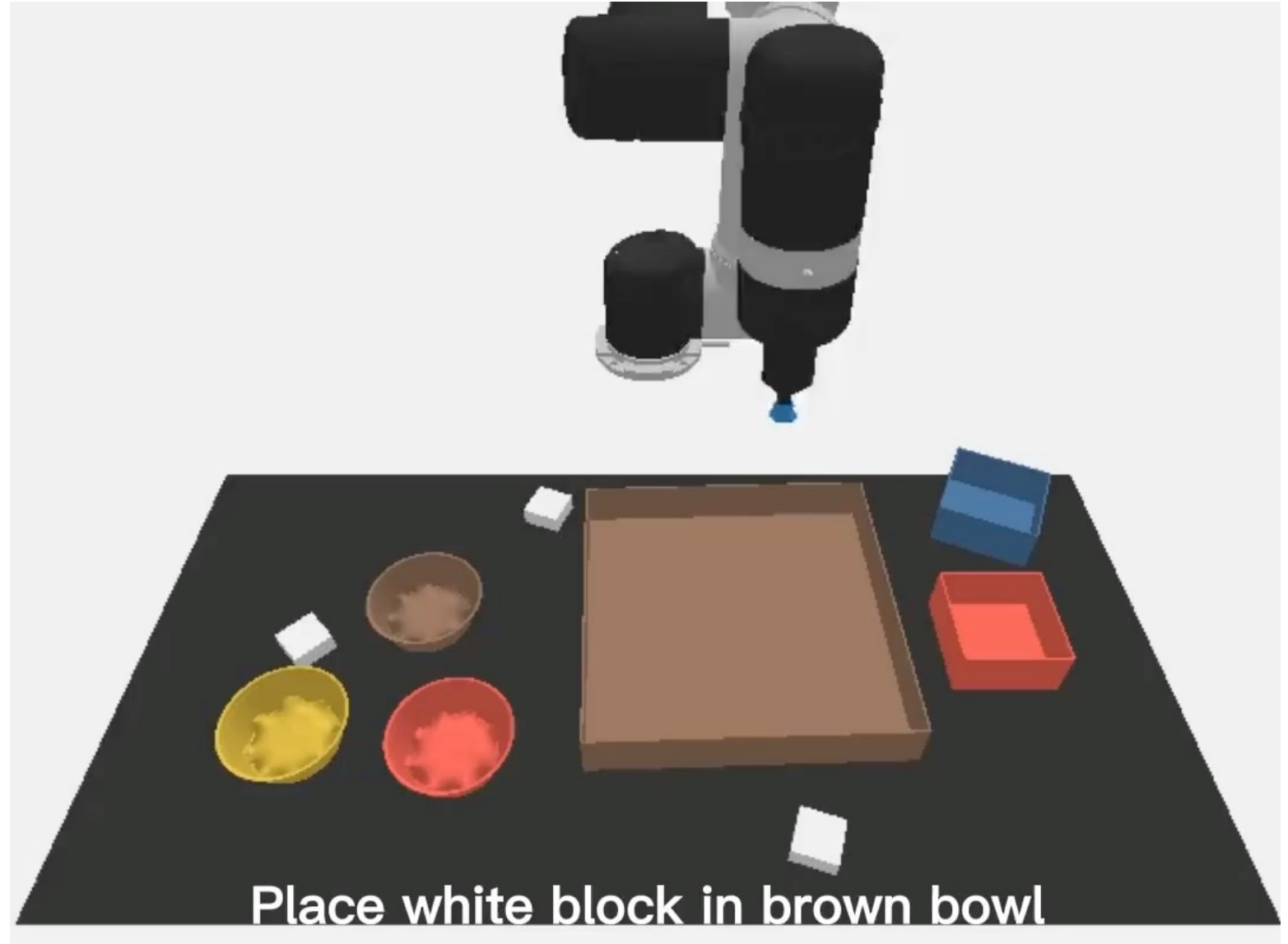
### Start Image



## Generated Language Plan

Place white block in a cyan bowl.
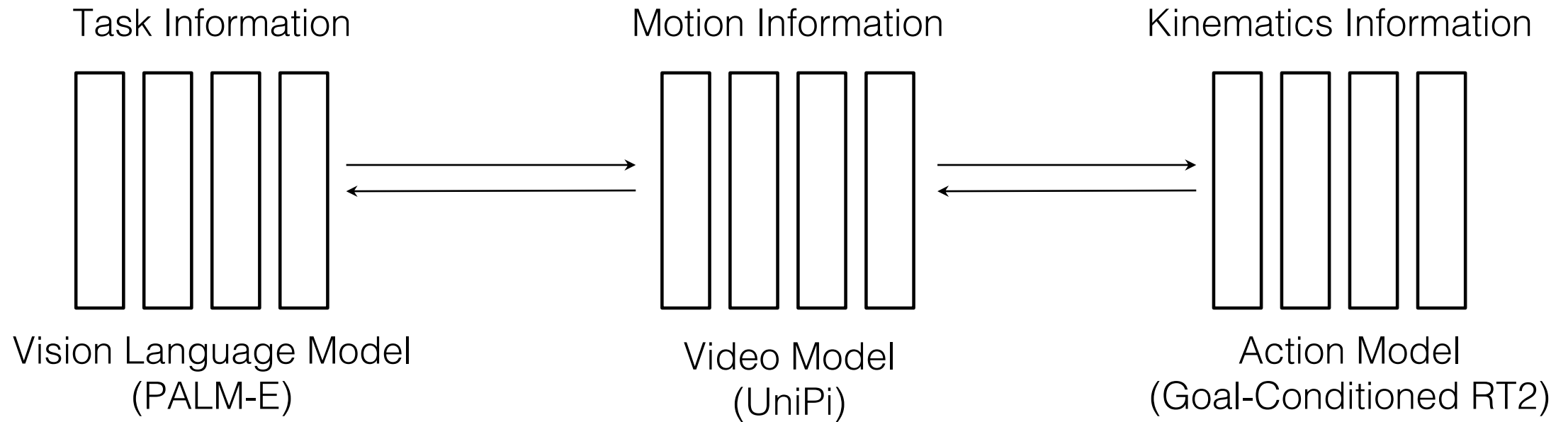
## Generated Video Plan

# Hierarchical Decision Making with Multimodal Models Execution

**Goal:** Stack red block on top of brown block and place yellow block to the left of the stack



Place white block in brown bowl

# Zero-Shot Planning and Execution with Foundation Models

Task Information

Motion Information

Kinematics Information

Vision Language Model
(PALM-E)

Video Model
(UniPi)

Action Model
(Goal-Conditioned RT2)

We can plan and execute unseen long horizon tasks on the real robot
without any explicit task training!

[1] Du et al. Video Language Planning. ICLR 2024

# Zero-Shot Planning and Execution with Foundation Models

**Goal:** Put the fruit into the top drawer.

# Overview

- Introduction to Multimodal Generative Models
- Compositional Generative Models
- Discussion

# Lecture 9: Generative AI, Part 1

## Yilun Du

CS 2281R: Mathematical & Engineering Principles for
Training Foundation Models