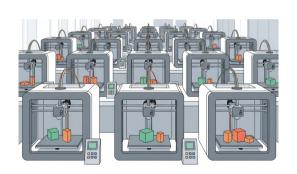
Speak to the Policy: Personalize MDP Policies with Language-based Priors

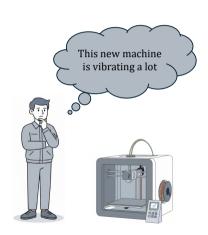
Qiyuan Chen and Raed Al Kontar

Industrial and Operations Engineering, University of Michigan, Ann Arbor

October 26, 2025

The Problem: How to "Speak" to Your Policy?

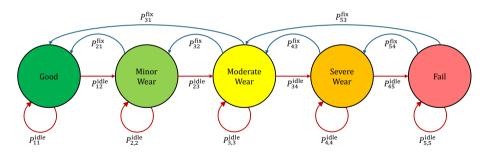




2 / 18

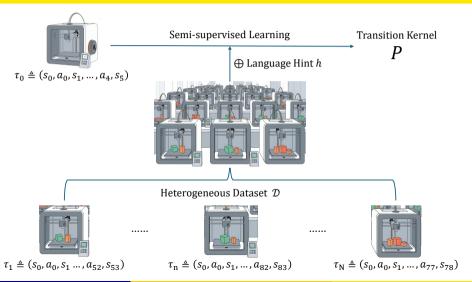
Background on MDP (An Example of Machine Maintenance)

- (Given) action $a \in A$, state $s \in S$, reward function: R(s, a)
- ullet (To Be Learned) transition dynamics: $P^a_{ii}=\mathbb{P}(s_{t+1}=j\mid s_t=i,a)$



Qiyuan Chen (UMich IOE)

Problem Formulation: Data Flow



Qiyuan Chen (UMich IOE)

Problem Formulation: The Inputs

We are given three ingredients:

- Past Data
 - The trajectory in the current environment: au_0
 - n historical trajectories from other environments: $\mathcal{D} = \{m{ au}_1, \dots, m{ au}_n\}$
- Human Language Hint
 - A text hint h describing the current environment behavior
- A Language Model
 - An LLM that can understand h

Goal

Learn a decision policy for the current environment with the help of the hint h the dataset $\mathcal D$

Overall Approach: A Two-Stage Pipeline

Our algorithm consists of two main stages:

Stage 1: Generate Priors via LLM

- Use the LLM and hint h to "judge" how likely trajectories lies in the main component
- ullet Train a model to produce a prior score s_k for each trajectory $oldsymbol{ au}_k$

Stage 2: Robust Kernel Estimation

- Use a semi-supervised, prior-weighted EM algorithm
- ullet This algorithm combines the labeled data $oldsymbol{ au}_0$ and the unlabeled data ${\cal D}$
- The priors generated in Stage 1 guide the EM algorithm to filter out outliers

Qiyuan Chen (UMich IOE) Speak to the Policy October 26, 2025 6/18

Step 1.1: LLM as a Judge

Gemini I am comparing two trajectories of a Markov Decision Process, which one is more likely given the hint of the environment? Background Knowledge: This is a machine maintenance problem. States of the MDP: 0 (machine is in good condition), 1 (minor wear) 2 (moderate wear) 3 (severe wear) 4 (fail) Actions of the MDP: M (maintenance), N (do nothing) Hint: "This new machine is vibrating a lot." Trajectory 1: (0.N.O.N.O.N.O.N.1.N.1.N.1.N.1.N.2.N.2) Trajectory 2: (0,N,0,N,1,N,1,N,2,N,3,N,4,M,0,N,1) Answer with only (A.B. or C) by choosing from the following options. A. Trajectory 1 B. Trajectory 2 C. Equally Likely / Cannot Determine



Collecting LLM's judgments

We input the following to the LLM:

- Background Knowledge
- Trajectory 1: τ_i
- Trajectory 2: τ_j
- Hint: h
- Prompt

The LLM outputs a preference:

- A. τ_i is more likely
- B. au_j is more likely
- C. Cannot Determine

Step 1.1: LLM as a Judge (cont.)

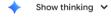
Gemini



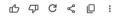
Hint: "This new machine is vibrating a lot."

Trajectory 1: (0,N,0,N,0,N,1,N,1,N,1,N,1,N,2,N,2)

Trajectory 2: (0,N,0,N,1,N,1,N,2,N,3,N,4,M,0,N,1)



В



Stage 1: Generating Language-based Priors

Stage 1: Generate Priors via LLM

- Step 1.1 collects a dataset of M pairwise preferences, $M = \mathcal{O}(n \log_2 n)$
- Step 1.2 trains a scoring function $f_{\theta}(\tau)$ that maps any trajectory τ to a scalar score

Bradlev-Terry Model

This model finds scores $f_{\theta}(\tau)$ such that the probability of τ_i being preferred to τ_i is a logistic function of their score difference:

$$\mathbb{P}(\boldsymbol{\tau}_i \succ \boldsymbol{\tau}_j) = \sigma(f_{\theta}(\boldsymbol{\tau}_i) - f_{\theta}(\boldsymbol{\tau}_j))$$

where $\sigma(z) = \frac{1}{1+e^{-z}}$ is the standard sigmoid function.

9 / 18

Step 1.2: Optimizing the Trajectory Score Model

We minimize the negative log-likelihood (given by B-T Model) of the M observed preferences:

$$\min_{ heta} \quad -rac{1}{M} \sum_{(oldsymbol{ au}_i \succ oldsymbol{ au}_j)} \log rac{\exp(f_{ heta}(oldsymbol{ au}_i))}{\exp(f_{ heta}(oldsymbol{ au}_i)) + \exp(f_{ heta}(oldsymbol{ au}_j))}$$

Result: Prior Scores

After training, we can compute a score for every unlabeled trajectory:

$$f_{\theta}(\boldsymbol{\tau}_k), \quad \forall k \in \{1, \ldots, n\}$$

This s_k is our prior belief that τ_k belongs to the main component.

Stage 2: Model Assumptions

We use a semi-supervised Expectation-Maximization (EM) algorithm.

Mixture Model

We assume each trajectory τ_k is drawn from a two-component mixture, determined by a latent variable $c_k \in \{\text{main, outlier}\}$.

- If $c_k = \text{main}$: τ_k is generated from P (our goal).
- If $c_k =$ outlier: τ_k is generated from a fixed uniform noise component.

Semi-Supervised

- For labeled data τ_0 , we fix $c_0 = \text{main}$.
- For unlabeled data τ_k , c_k is a latent variable.

Model Assumptions (Visual)



Components:

- Main component (black)
- Observed samples (red)
- Background noise (gray)

Correspondence:

- τ_0 : Red
- \mathcal{D} : Black+ Gray

Stage 2: EM Algorithm - Initialization

Step 1: Calculate Priors (for $\tau_k \in \mathcal{D}$)

Convert scores to prior probabilities:

$$\pi_k \triangleq \mathbb{P}(c_k = \mathsf{main}) \approx \sigma(f_\theta(\boldsymbol{\tau}_k) - \kappa),$$

where κ is a hyperparameter that controls the background noise level.

* A simple choice of κ can be $\kappa = \frac{1}{n} \sum_{k} f_{\theta}(\tau_{k})$.

Step 2: Initialize Main Kernel P

Set initial estimate using *only* the labeled data τ_0 (a is suppressed):

$$P_{ij} = \frac{N_{ij}^k}{\sum_{l=1}^S N_{il}^k}$$

Stage 2: EM Algorithm - The Loop

We repeat the E-Step and M-Step until the kernel P converges.

E-Step (Expectation)

Calculate the posterior probability (responsibility) r_k that each trajectory τ_k belongs to the main component, given the current kernel P.

M-Step (Maximization)

Update the kernel P by calculating the weighted MLE, using the responsibilities r_k as weights.

Stage 2: EM Algorithm - E-Step

For Labeled Data (au_0)

The responsibility $r_0 = \mathbb{P}(c_0 = \mathsf{main} | \boldsymbol{ au}_0, P)$ is fixed to our certain belief: $r_0 = 1$

For Unlabeled Data ($au_k \in \mathcal{D}$)

Find the posterior responsibility using Bayes' rule:

$$r_k = \mathbb{P}(c_k = \mathsf{main}|\boldsymbol{\tau}_k, P) = \frac{\overbrace{\boldsymbol{\pi}_k}^{\mathsf{Langauge Prior}} \underbrace{\mathbb{P}(\boldsymbol{\tau}_k|c_k = \mathsf{main}, P)}^{\mathsf{Likelihood (Main)}}}{\frac{\mathbf{T}_k}{\mathbf{T}_k} \cdot \mathbb{P}(\boldsymbol{\tau}_k|c_k = \mathsf{main}, P) + (1 - \pi_k) \cdot \mathbb{P}(\boldsymbol{\tau}_k|c_k = \mathsf{noise}, P)}$$

$$\mathbb{P}(\boldsymbol{\tau}_k|c_k = \mathsf{main}, P) = \prod_{i=1}^S \prod_{j=1}^S (P_{ij})^{N_{ij}^{(k)}}, \quad \mathbb{P}(\boldsymbol{\tau}_k|c_k = \mathsf{noise}, P) = \left(\frac{1}{S}\right)^{|\boldsymbol{\tau}_k|}$$

15 / 18

Stage 2: EM Algorithm - M-Step

We update P by calculating the weighted MLE

- ullet "Hard" counts from labeled data au_0
- ullet "Soft" (responsibility-weighted) counts from unlabeled data ${\cal D}$

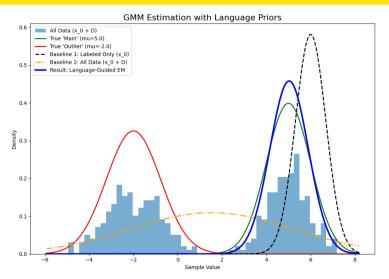
Update Rule

For each transition pair $i \rightarrow j$:

$$P_{ij} = \frac{N_{ij}^{(0)} + \left(\sum_{k \in \mathcal{D}} r_k \cdot N_{ij}^{(k)}\right)}{\sum_{l=1}^{S} \left[N_{il}^{(0)} + \left(\sum_{k \in \mathcal{D}} r_k \cdot N_{il}^{(k)}\right)\right]}$$

Repeat E-Step and M-Step until P converges.

Visualizations on a Gaussian Mixture Model



Toy Example:

- Demo with 1-d GMM
- Hint: "values are larger"
- $\kappa = \frac{1}{n} \sum_{k} f_{\theta}(\boldsymbol{\tau}_{k})$

Highlights:

- Beyond MDP
- Robust Recovery
- Non-uniform noise

17 / 18

Qiyuan Chen (UMich IOE) Speak to the Policy October 26, 2025

Thank You for Listening!

Speak to the Policy

- Stage 1 An LLM judges trajectories based on a human instruction h. A choice model is trained on these judgments to create prior scores.
- **Stage 2** The language prior guides the semi-supervised learning of the transition kernels. The model is optimized by an EM algorithm.



Qiyuan-Chen.com