

CS 343

Review and Conclusion

Prof. Yuke Zhu

The University of Texas at Austin



Announcements

- Please fill out the course survey
 - Feedback to both instructor and TAs
 - Positive and negative points are useful
 - Post on Piazza your completion screenshot (in a private post) as a form of participation!
- Capture the Flag contest results!

CTF Contest

15 teams participated, 10 qualified

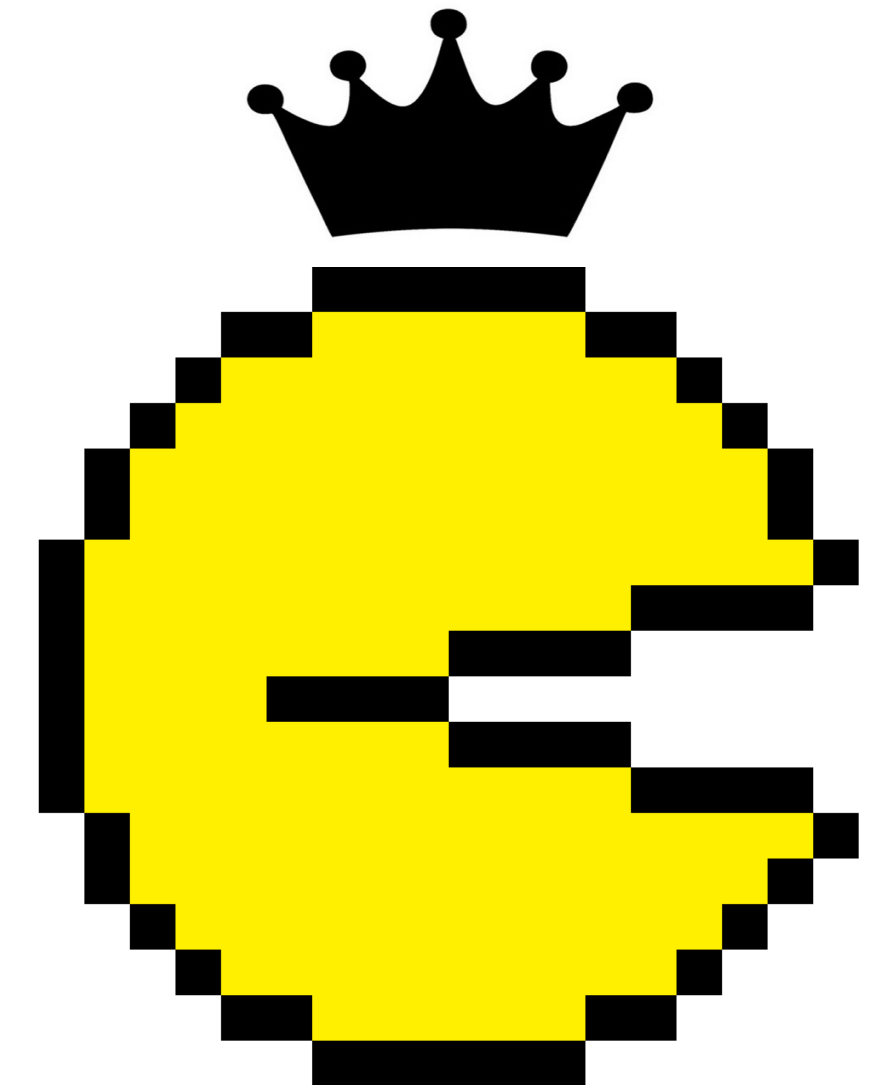
Rankings

1st place: Joseph Stanley, Ritvik Renikunta

2nd place: Adit Pareek, Eylam Tagor

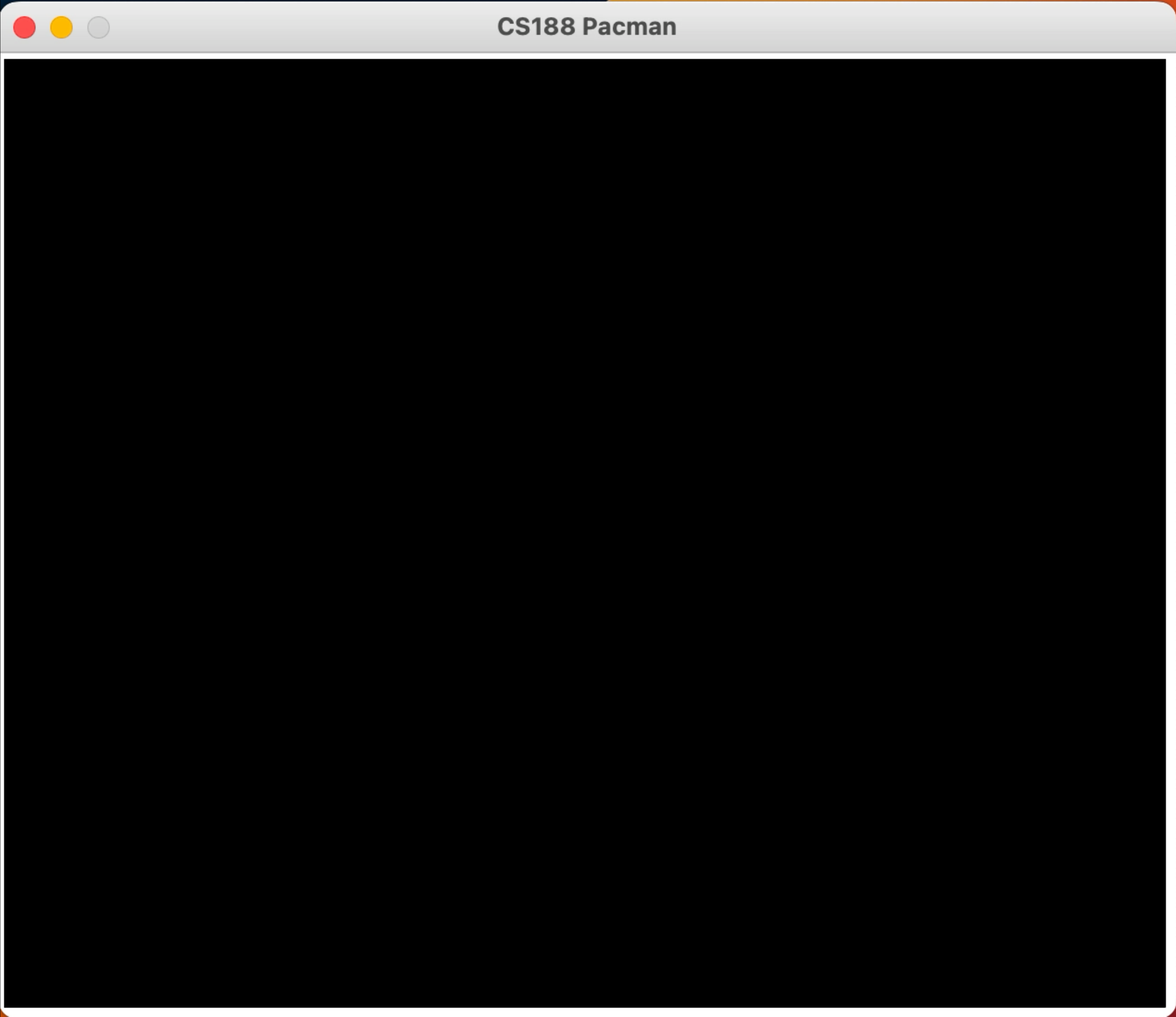
3rd place: Arik Rundquist (tie), John Li (tie)

Congratulations!



1st place

Joseph Stanley,
Ritvik Renikunta



2nd place

Adit Pareek,
Eylam Tagor

Overview of AI Topics

Search / Planning

Uninformed Search

A* Search

CSPs

Local Search

Minimax

Expectimax

MDPs

Machine Learning

Reinforcement Learning

Probability Theory

Bayes Nets

HMMs

Particle Filters

Decision Diagrams

Naive Bayes

Perceptrons

Neural Networks

Kernels

Clustering

VPI

Overview of Machine Learning

Supervised Learning

Discriminative Models

Perceptrons

Neural Networks

Generative Models

Bayes Nets

Naive Bayes

HMMs

Reinforcement Learning

MDPs

Value Iteration

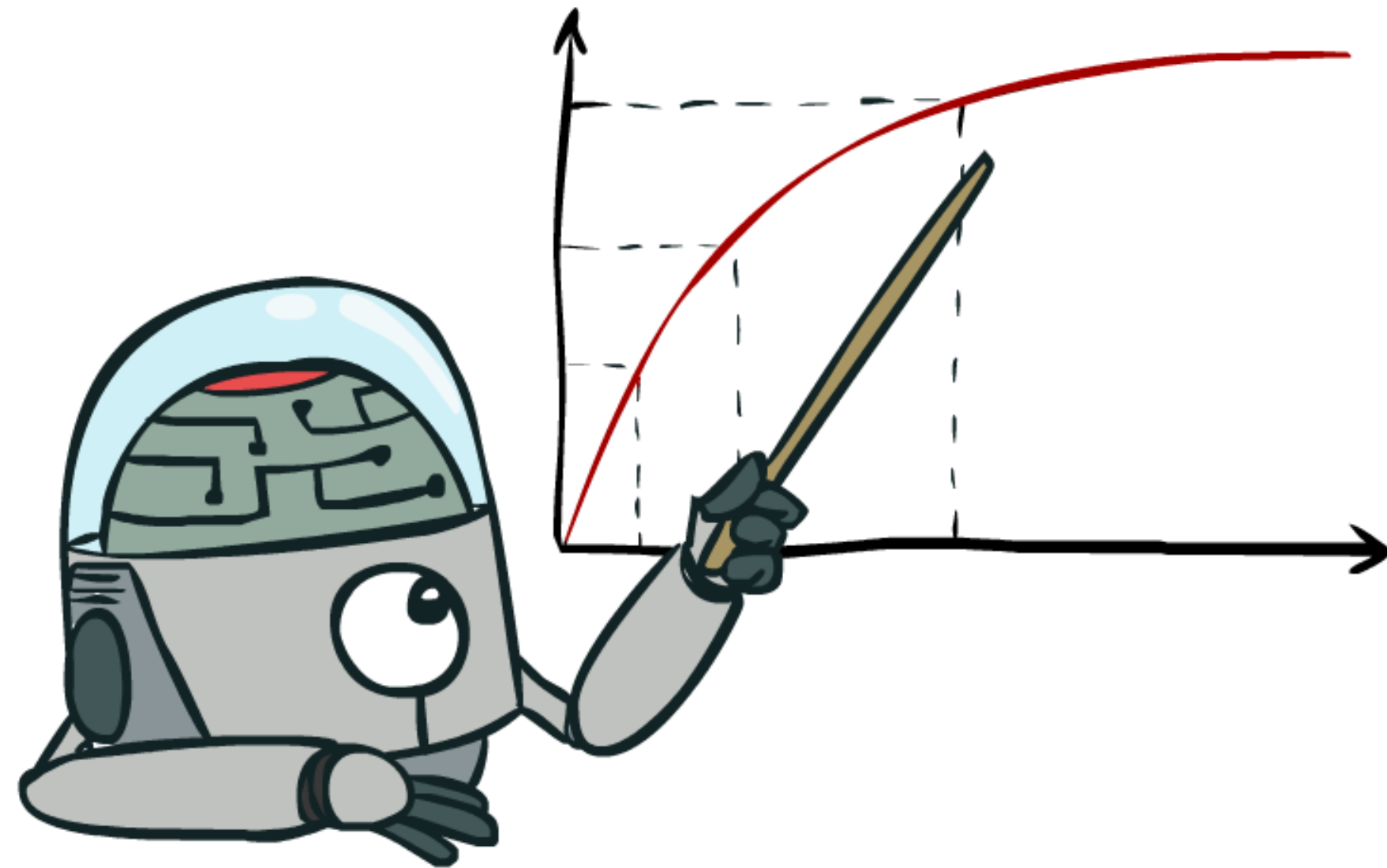
Policy Iteration

Q Learning

Unsupervised Learning

K-Means Clustering

Maximize Your Expected Utility



How Do AI Systems Maximize Utility?

Constraint satisfaction: searching intelligently for legal solutions

8			4		6			7
						4		
	1					6	5	
5		9		3		7	8	
				7				
	4	8		2		1		3
	5	2					9	
		1						
3			9		2			5

Example: Sudoku

Utility: Does the solution satisfy the rules / constraints?

Assumptions: We can write down the rules / constraints of the problem

How Do AI Systems Maximize Utility?

Planning: reasoning with models



Example: Robot navigation

Utility: Path length, collisions, surfaces, energy, social factors

Assumptions: We have a model of the world and the effects of the agent's actions

How Do AI Systems Maximize Utility?

Supervised Learning: learning from labeled examples



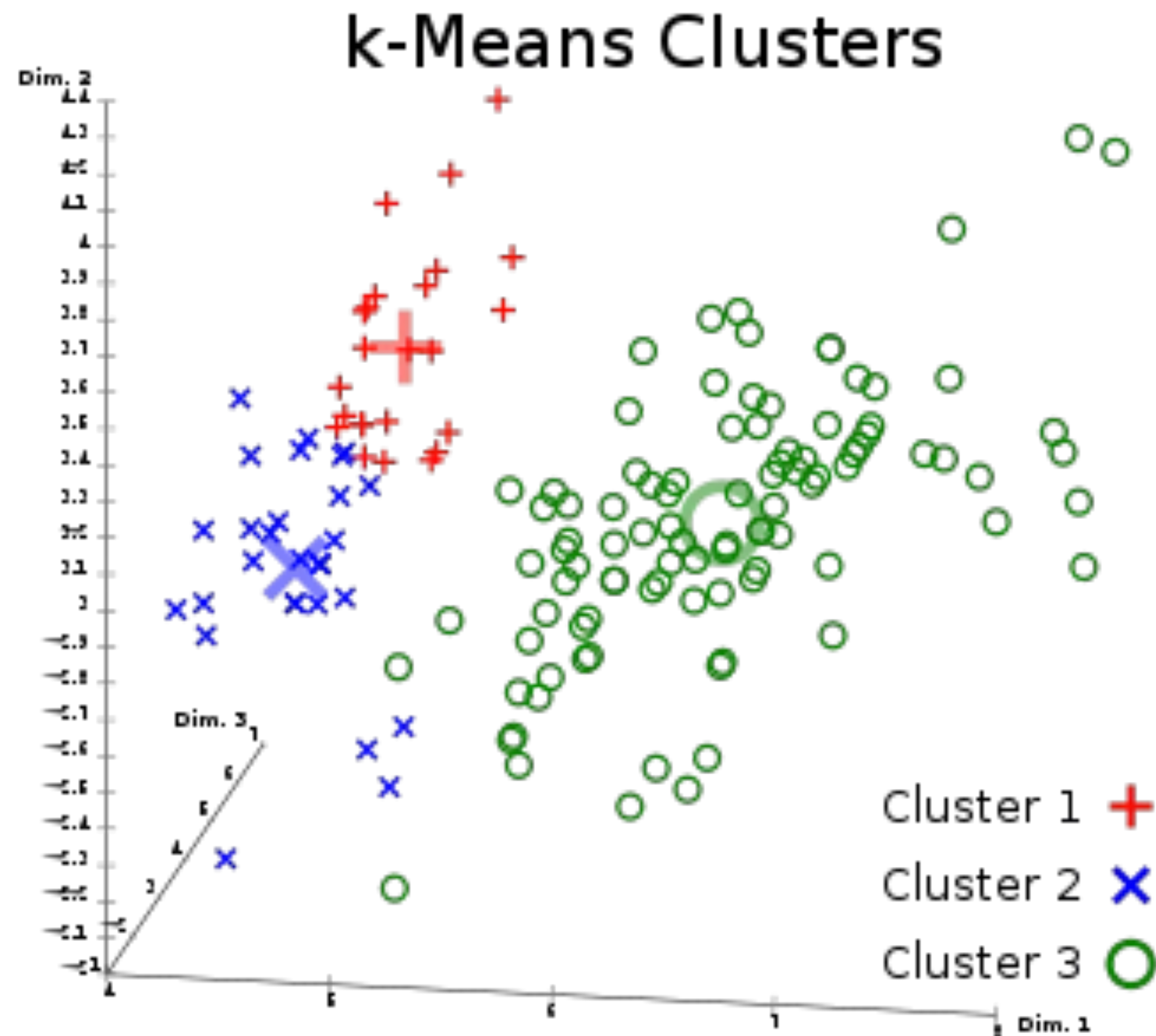
Example: Image classification

Utility: Classification accuracy on images not seen during training

Assumptions: We have access to a (usually large) labeled data set

How Do AI Systems Maximize Utility?

Unsupervised Learning: discovering patterns in unlabeled data



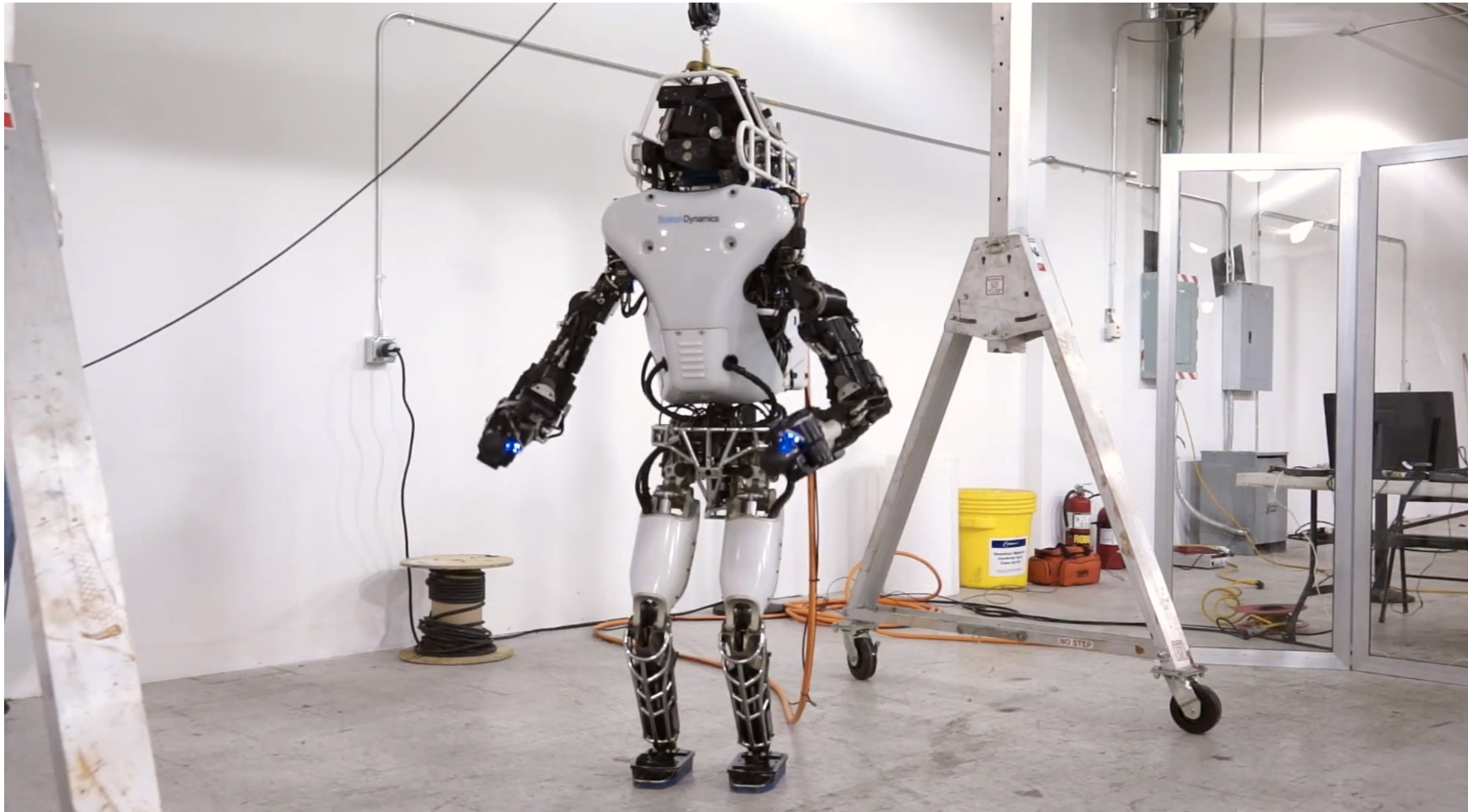
Example: Clustering species

Utility: “Best” explanation of data

Assumptions: Data points that should be clustered together are “close” together

How Do AI Systems Maximize Utility?

Reinforcement Learning: learning from experience



Example: Robot walking

Utility: Time until fall, speed, energy efficiency

Assumptions: We can “reward” and “punish” good and bad performance, but don’t know what the correct action at each step should be

Properties of task environment

- Single-agent vs. multi-agent
- Deterministic vs. stochastic
- Fully observable vs. partially observable
- Episodic vs. sequential
- Static vs. dynamic
- Discrete vs. continuous
- Known vs. unknown

Single agent vs. multi-agent

- Not multi-agent if other agents can be considered part of the environment
- Only considered to be multi-agent if the agents are maximizing a performance metric that depends on other agents' behavior
- Single agent example: Pacman with randomly moving ghosts
- Multi-agent example: Pacman with ghosts that use a planner to follow him

Single / Multi Agent

Single

Uninformed Search

A* Search

Local Search

CSPs

Multi

Minimax

Expectimax

MDPs

Reinforcement Learning

Deterministic vs. stochastic

- Deterministic: next state of environment is completely determined by the current state and the action executed by the agent
- Stochastic: actions have probabilistic outcomes
- Strongly related to partial observability — most apparent stochasticity results from partial observation of a deterministic system
- Example: Coin flip

Determinism

Deterministic

Uninformed Search

A* Search

Local Search

CSPs

Minimax

Stochastic

Expectimax

MDPs

Reinforcement Learning

Decision Diagrams

Fully observable vs. partially observable

- Fully observable: agent's sensors give it access to complete state of the environment at all times
- Can be partially observable due to noisy and inaccurate sensors, or because parts of the state are simply missing from the sensor data
- Example: Perfect GPS vs noisy pose estimation
- Example: IKEA assembly while blindfolded

Almost everything in the real world is partially observable

Observability

Fully Observable

Uninformed Search

A* Search

Local Search

CSPs

Minimax

Expectimax

MDPs

Reinforcement Learning

Partially Observable

Bayes Nets

HMMs

Decision Diagrams

Known vs. unknown

- Agent's state of knowledge about the "rules of the game" / "laws of physics"
- Known environment: the outcomes for all actions are given
- Unknown: agent has to learn how it works to make good decisions
- Possible to be partially observable but known (solitaire)
- Possible to be fully observable but unknown (video game)

Model of the World

Known

Uninformed Search

A* Search

Local Search

CSPs

Classic Planning

Minimax

Expectimax

MDPs

Value Iteration

Decision Diagrams

Unknown

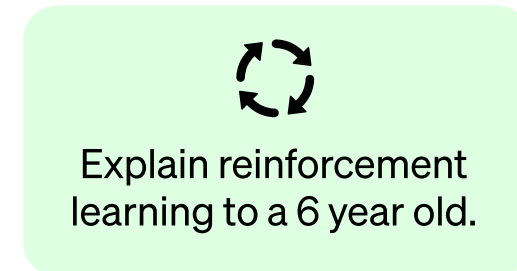
Q Learning

Learning parameters
of Bayes Net

Step 1

Collect demonstration data and train a supervised policy.

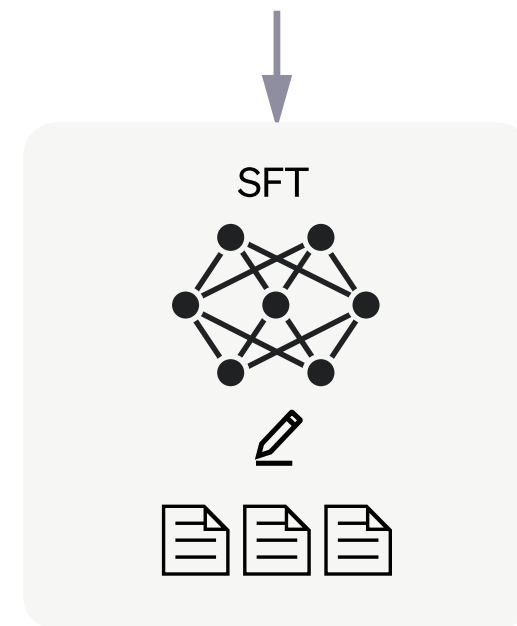
A prompt is sampled from our prompt dataset.



A labeler demonstrates the desired output behavior.



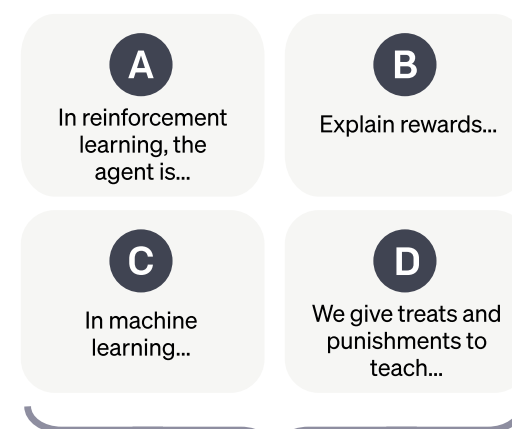
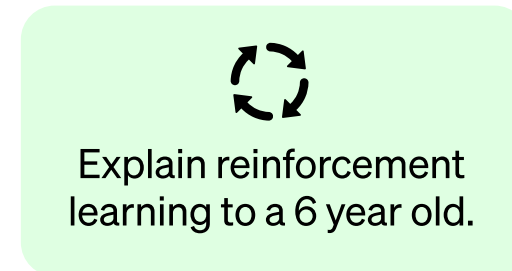
This data is used to fine-tune GPT-3.5 with supervised learning.



Step 2

Collect comparison data and train a reward model.

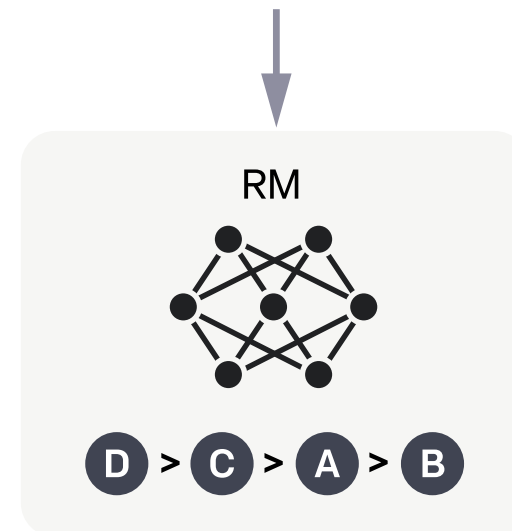
A prompt and several model outputs are sampled.



A labeler ranks the outputs from best to worst.



This data is used to train our reward model.



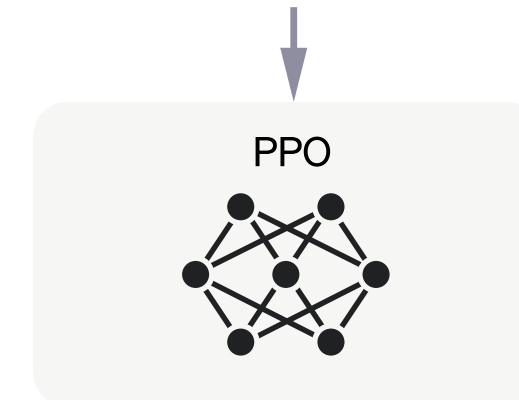
Step 3

Optimize a policy against the reward model using the PPO reinforcement learning algorithm.

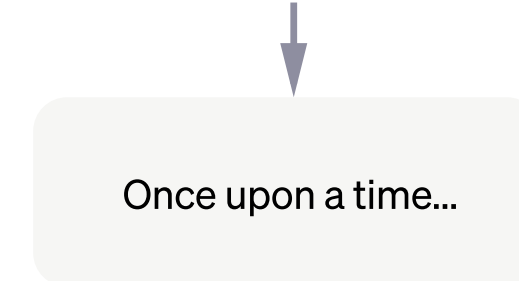
A new prompt is sampled from the dataset.



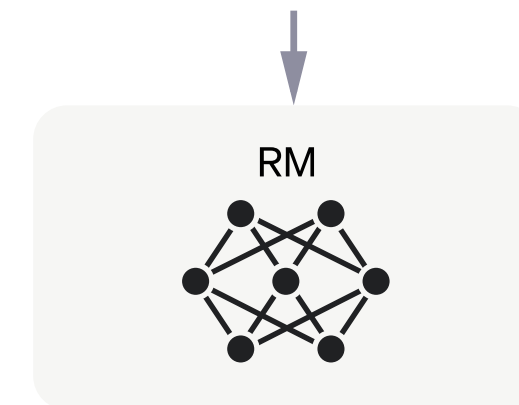
The PPO model is initialized from the supervised policy.



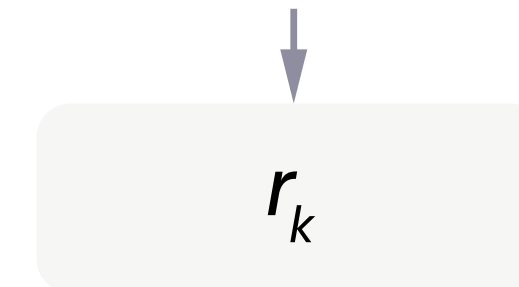
The policy generates an output.



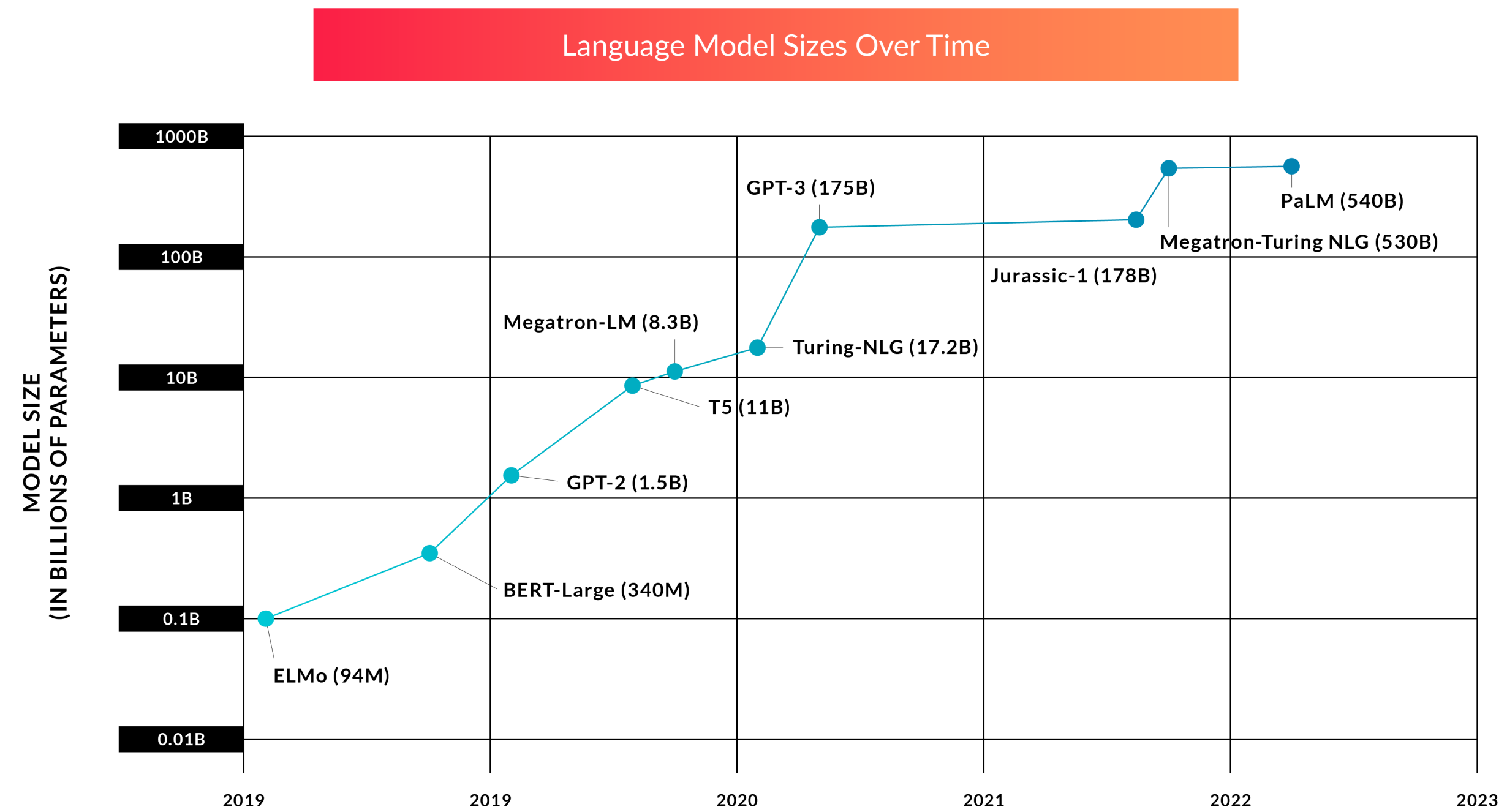
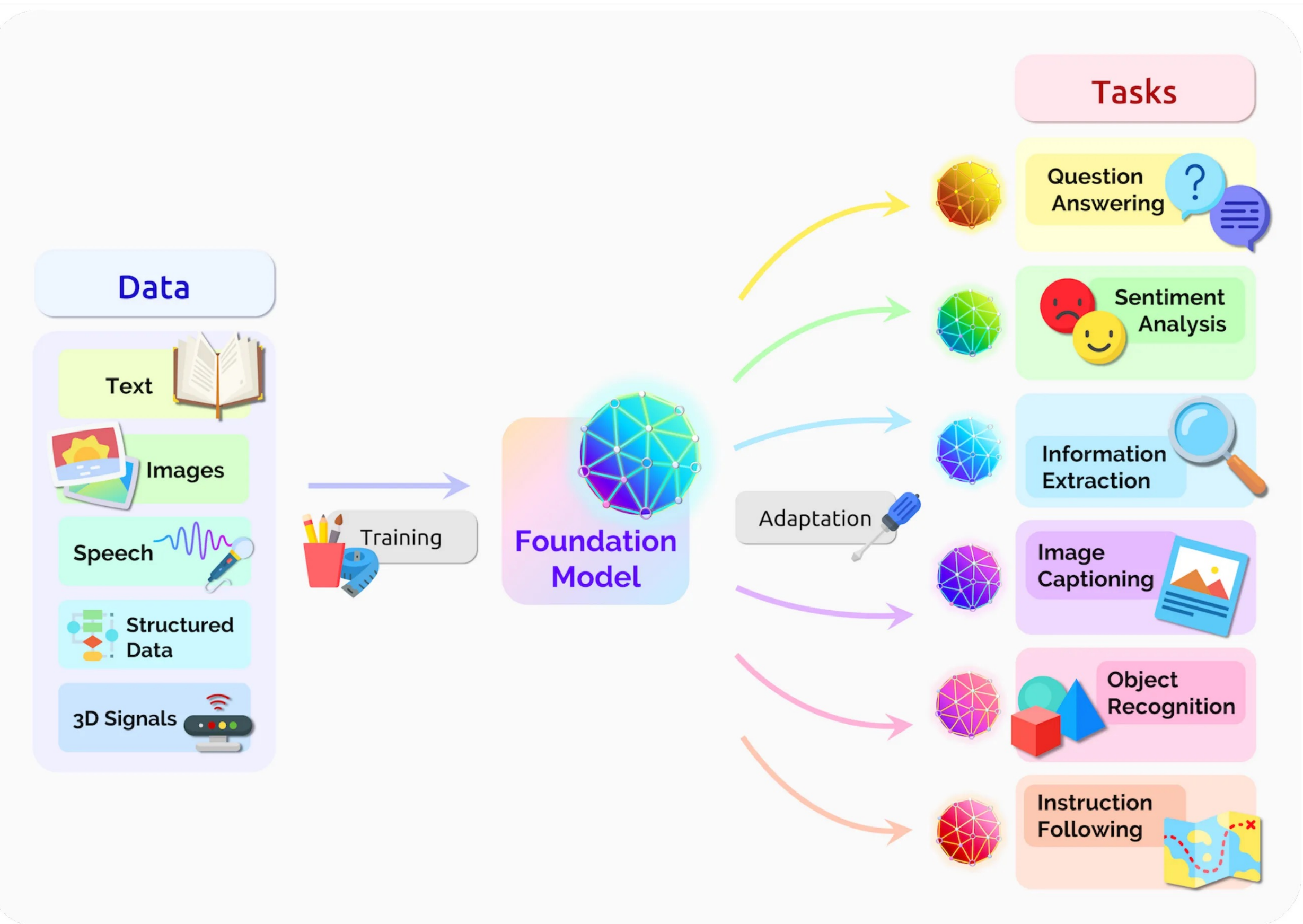
The reward model calculates a reward for the output.



The reward is used to update the policy using PPO.



How ChatGPT is Built



Era of Big Models

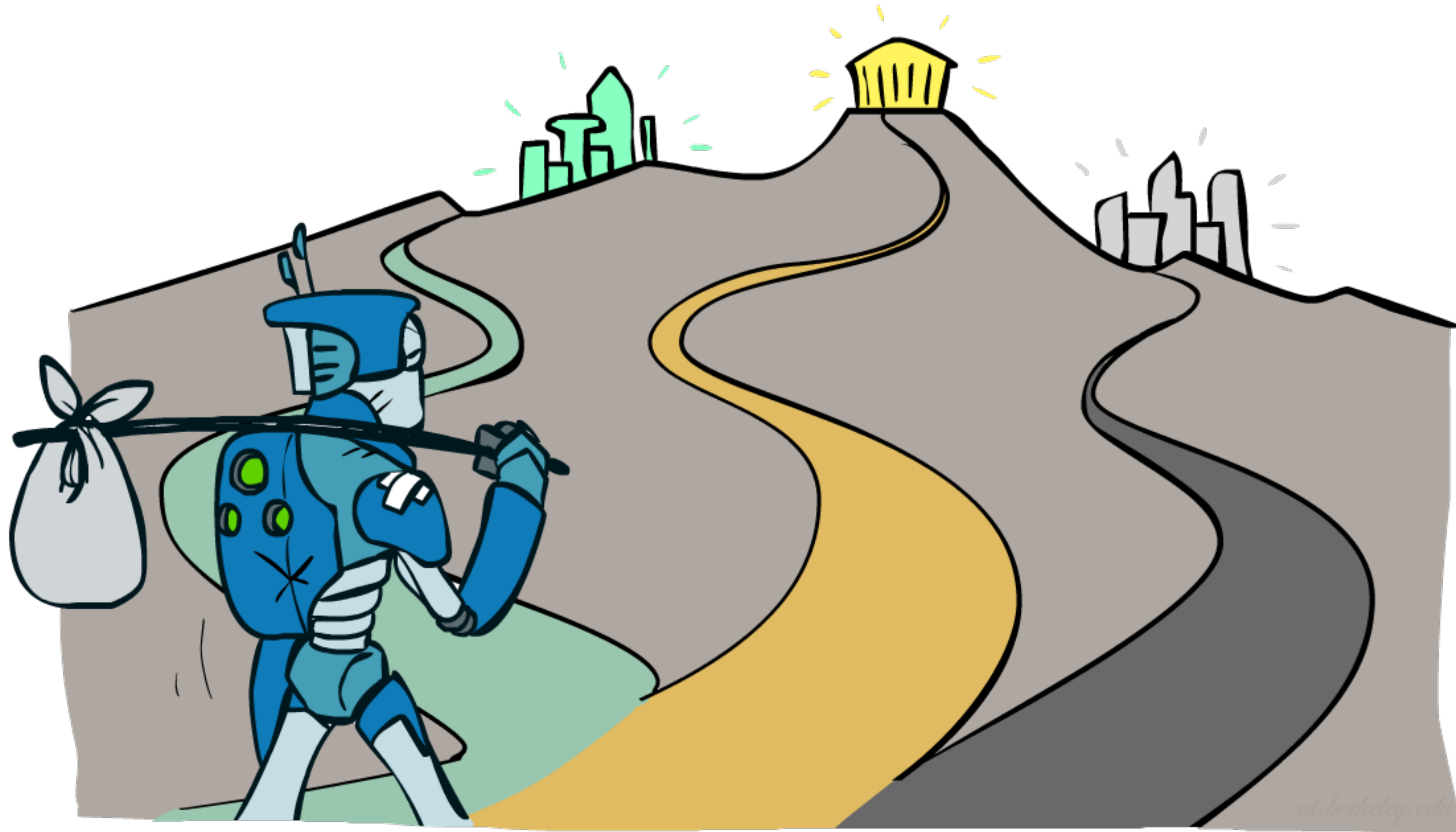


Workspace Image

2x

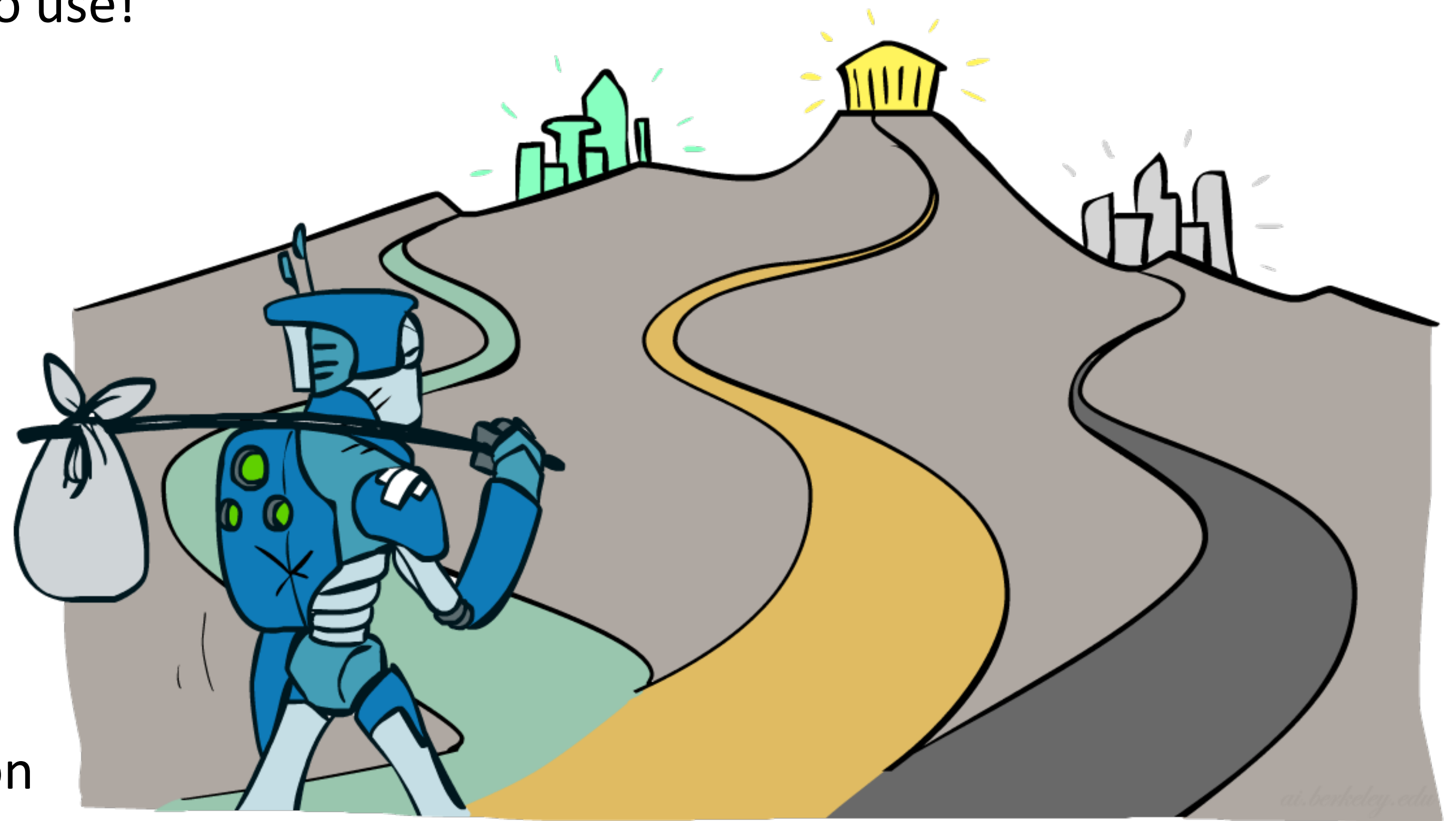
[Credit: Steve Han and Mingyo Seo]

Where to Go Next?



Where to Go Next?

- Congratulations, you've seen the basics of modern AI
 - ... and done some amazing work putting it to use!
- How to continue:
 - CS 395T Visual Recognition
 - CS 391R Robot Learning
 - ECE 382V Human Robot Interaction
 - CS 388 Natural Language Processing
 - CS 391L Machine Learning
 - CS 393R Autonomous Robots
 - CS 342 Neural Networks
 - EE 381V Advanced Topics in Computer Vision
 - CS 394R Reinforcement Learning: Theory and Practice
 - ... and more; ask if you're interested



Final Remarks

- We have come a long way! Thank you!
- We are very proud that you have made it to the end of this demanding course!
- We are impressed by your ingenuity and critical thinking in the in-class discussions, Piazza posts, projects, and assignments!
- Thanks to Zhenyu and Pranav for handling the course logistics.
- If this course helps you kickstart your future endeavors in AI, please email us and let us know!

Thank you

Friday 4/28 8 – 10am GEA 105

1 page (front and back) of notes

Closed book

That's it!

I had a great time teaching this course and
I hope you all enjoyed it as well

Have a great summer!