

DS-GA 3001 007 | Lecture 6

Reinforcement Learning

Jeremy Curuksu, PhD

NYU Center for Data Science

jeremy.cur@nyu.edu

March 2, 2023

DS-GA 3001 RL Curriculum

Reinforcement Learning:

- ▶ Introduction to Reinforcement Learning
- ▶ Multi-armed Bandit
- ▶ Dynamic Programming on Markov Decision Process
- ▶ Model-free Reinforcement Learning
- ▶ Value Function Approximation (Deep RL)
- ▶ **Examples of Industrial Applications and Project Q&A**
- ▶ Policy Function Approximation (Actor-Critic)
- ▶ Planning from a Model of the Environment
- ▶ Advanced Topics and Development Platforms

Reinforcement Learning

Last week: Value Function Approximation

- ▶ Categories of Functions in Reinforcement Learning
- ▶ Approximation of Value Functions
- ▶ Deep Reinforcement Learning

Today: Examples of Industrial Applications

- ▶ A Tour of 10 Awesome Applications of RL
- ▶ Project Q&A

Robotics

Teach a Robot to ...



Source: DeepMind (2022)

Autonomous Driving

Learn to Drive Like a Human

- **Goal:** Drive vehicle on a circuit to destination without leaving the road



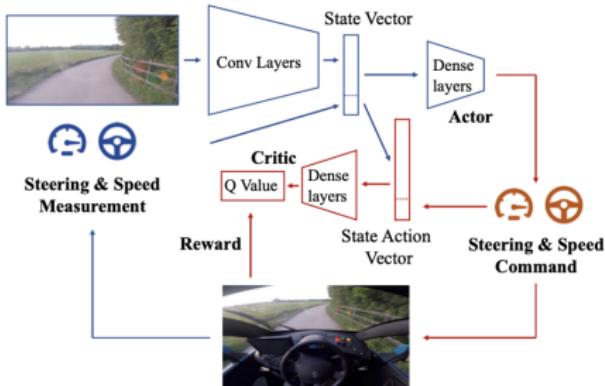
Environment:

Virtual and physical driving lanes

Reward Function:

- Forward Speed
- Termination upon *infraction of traffic rules* by safety driver

DS-GA 3001 007 | Lecture 6



State:

- Pixels of front camera encoded by CNN (single monocular image)
- Vehicle speed and steering angle

Actions:

2 actions: Speed, steering angle

Learn to Drive Like a Human

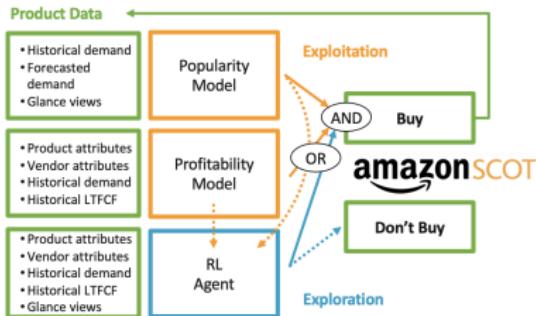


Source: Wayve (2019)

Amazon Inventory Management

Manage Amazon Retail Inventory

- **Goal:** Select products to show as "shipped and sold by Amazon" on the Amazon.com website to maximize customer experience and profitability



Offline Validation: All auto parts and medical supply products in catalog as of 8/11/18.
Cash flow and sales collected for year after 8/11/18.

Statistics	ASINs selected to buy	ASINs blocked from buying
ASINs count (in MM)	0.16 (2.7%)	5.55 (97.3%)
Cash flow (in MM euros)	152.35	-19.04
Sales (in MM)	28.06 (91.7%)	2.57 (8.3%)

Online A/B testing: Q4 2019, 30M products, 90% treatment, 10% control

EU LAB	Treatment Effect per ASIN per week	Confidence Interval	p-value	Annualized Impact
CP (Euros)	0.0103	[0.002, 0.019]	0.02	€2.45 MM
Sales (Euros)	0.021	[-0.061, 0.103]	0.10	€4.68 MM
Cash flow (Euros)	0.1123	[-0.311, 0.536]	0.54	€25.03 MM
Out of stock (bps)	-74	[-100, -50]	0.00	-74 bps

Environment:

Cash flow and popularity, $\Delta t = 3$ months

Reward Function:

$$\text{Cash flow} = \sum \left(\begin{array}{l} \text{short term profit} + \text{long term value} \\ -\text{cost of capital} - \text{asset depreciation} \end{array} \right)$$

State:

Product-, brand- and vendor-level statistics, historical sales, historical cash flow, glance views

Actions:

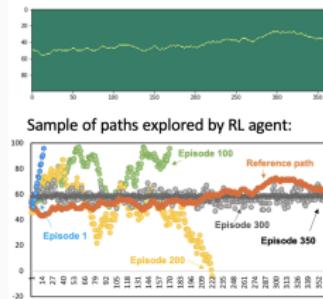
1 action: (block buy, buy at least 1 unit)

Seismic Mapping to Identify Natural Oil & Gas Reserves

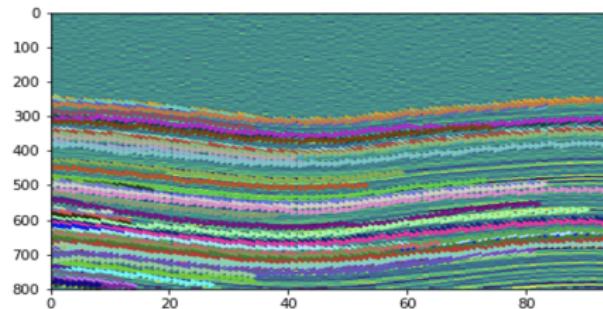
Seismic Pathfinder

- **Goal:** Screen cross-sections under the earth surface to identify the nature and geometry of individual seismic layers, to reduce exploration costs

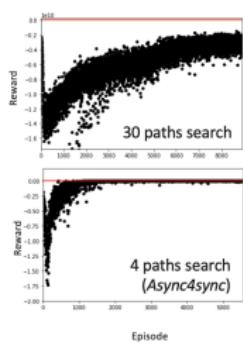
Pathfinder for 1 path:



1 cross-sectional (cylinder) map:



Accumulated reward:



Environment:

Underground cylinder cross-sections

Reward Function:

$$-\beta_1 \sum |z_t^i - z_{ref}^i| - \beta_2 \sum |z_{end}^i - z_0^i|$$

State:

Indirect seismic measurements at z_t and $(z-3, z-2, z-1, z, z+1, z+2, z+3)_{t+1}$

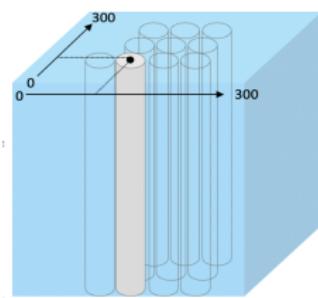
Actions:

k actions: (up, down) for each path

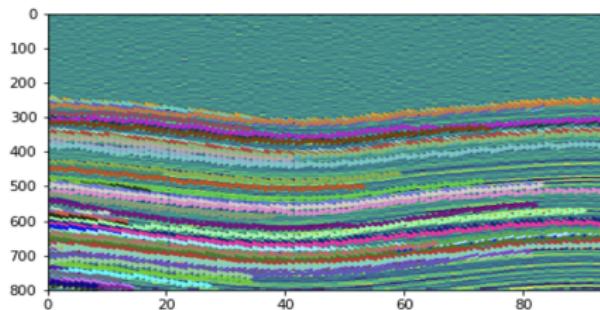
Seismic Pathfinder

- **Goal:** Screen cross-sections under the earth surface to identify the nature and geometry of individual seismic layers, to reduce exploration costs

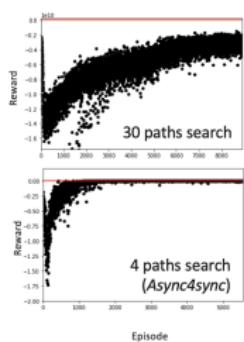
3D seismic cube:



1 cross-sectional (cylinder) map:



Accumulated reward:



Environment:

Underground cylinder cross-sections

Reward Function:

$$-\beta_1 \sum |z_t^i - z_{ref}^i| - \beta_2 \sum |z_{end}^i - z_0^i|$$

State:

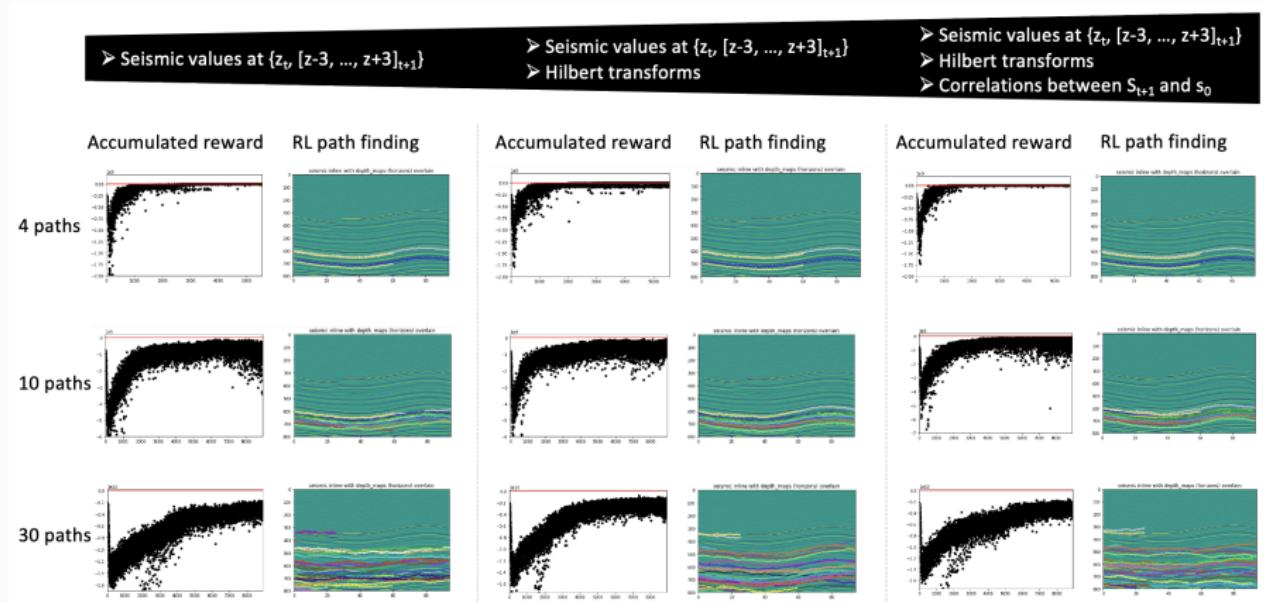
Indirect seismic measurements at z_t and $(z-3, z-2, z-1, z, z+1, z+2, z+3)_{t+1}$

Actions:

k actions: (up, down) for each path

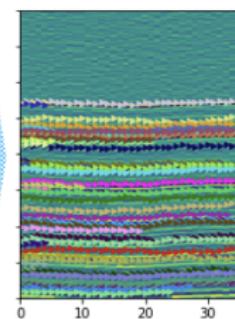
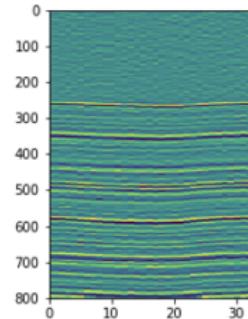
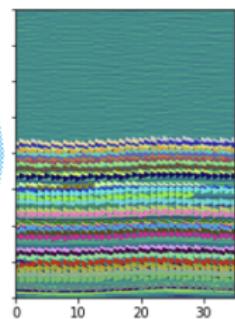
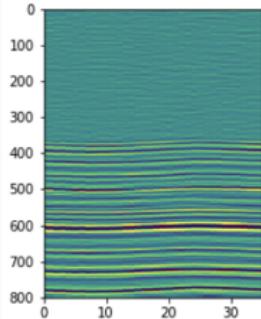
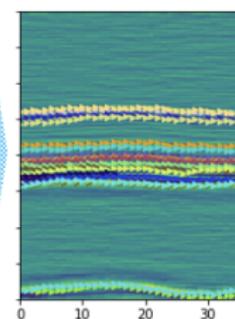
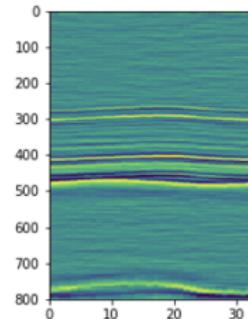
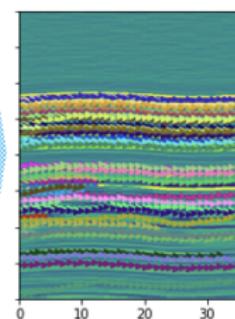
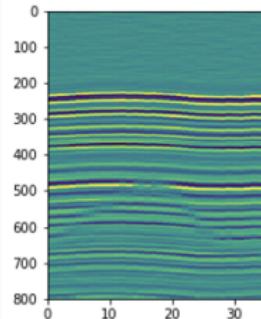
Seismic Pathfinder

- ▶ **Analysis of Results:** Benchmarks of synchronous n -path search RL on state complexity and number of paths



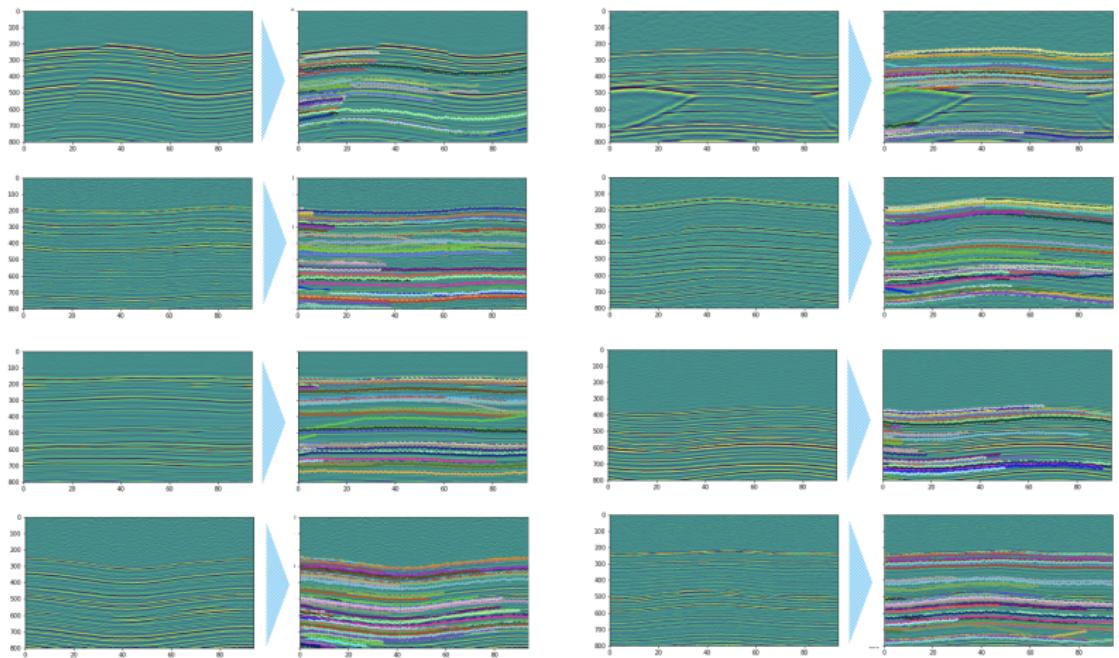
Seismic Pathfinder

- ▶ **Analysis of Results:** Generalization of pre-trained Async4sync DRL agent on arbitrary cubes and cylinders with **radius = 35 steps**



Seismic Pathfinder

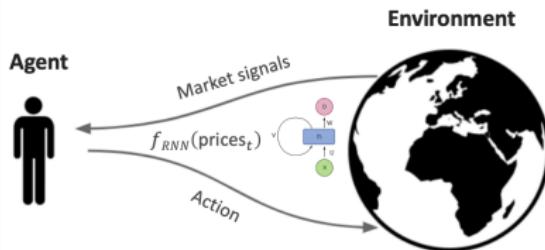
- ▶ **Analysis of Results:** Generalization of pre-trained Async4sync DRL agent on arbitrary cubes and cylinders with **radius = 90 steps**



Algorithmic Trading

Manage a Stock Portfolio

- **Goal:** Identify trading strategy in portfolio of k stocks to maximize profit



RNN-based RL outperforms lag-based RL at trading stocks

Mean return \$4,900 vs. \$2,200; Maximum return \$7,500 vs. \$2,800



Environment:

7 years of stock prices and news headlines, $\Delta t = 1$ day

State:

Vector of k stock prices for last n days
encoded by RNN

Reward Function:

$$r_t = r_t^0 + r_t^{risk} + r_t^{fee} = \sum_k \underbrace{\frac{(\text{prices}_{t+1} - \text{prices}_t)}{\text{prices}_t}}_{\text{Portfolio Value Change}} x_t - \lambda \sigma_t^2(r_t^0) - \kappa_t^T x_t$$

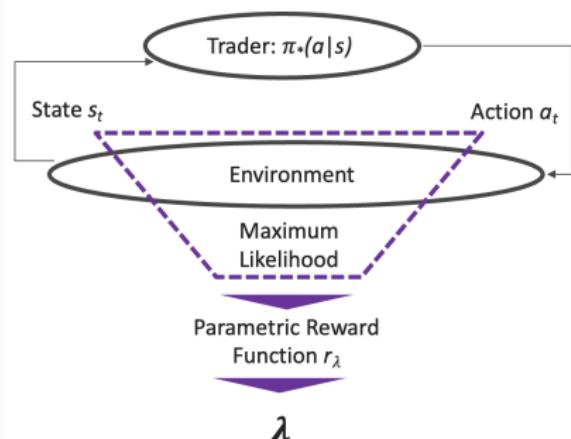
Risk Penalty Transaction Cost

Actions:

k actions: (\$buy, \$sell, sit) $_k$

Interpret Financial Trading Behavior

- ▶ **Inverse Reinforcement Learning:** Identify trader attributes, such as a level of risk aversion, by observing its behaviors
- ▶ **Estimate parameters of a reward function** that fit observed trajectories under a given policy in historical or simulated experience
- ▶ **Example:** Compute the risk aversion parameter λ of a successful trader



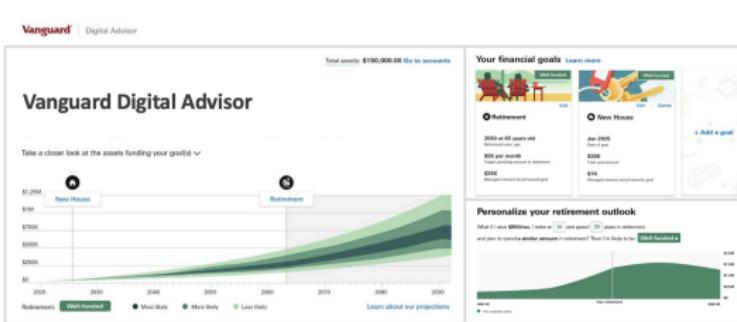
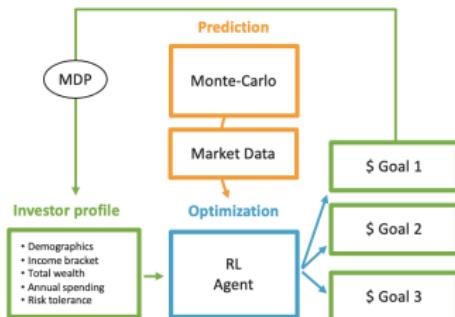
Reverse Engineer Strategy from Trader (= proprietary trading events) using Inverse RL:

1. Choose a parametric form of reward function
2. Estimate its parameters (MLE) from observed behavior in past trading events until convergence
3. Apply RL under the estimated reward function on an arbitrary stock portfolio, to identify an optimal policy (trading strategy) for this stock portfolio

Asset Allocation and Wealth Management

Manage Long-Term Financial Goals

- **Goal:** Determine optimal asset allocation strategy to meet multiple long-term financial goals, while also being successful in retirement



Environment:

50,000 investors profiles

$$\mathbb{E}(s' | s, a) + w(\sigma^2) \text{ and } \mathbb{E}(p_{100}), \Delta t = 1 \text{ year}$$

Reward Function:

$$r_t = r_t^{work} + r_t^{retire} = -\beta_1 \sum |g_t^i - g_T^i| + \beta_2 p_{100}$$

State:

Investor profile (age, location, income, wealth, spending, risk tolerance), \$ contributed for each goal (g_t^i), time left

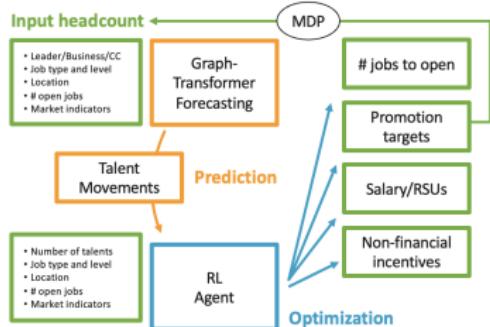
Actions:

k actions: \$ contribution for each goal

Workforce Management

Manage a Large Talent Workforce

- **Goal:** Pull talent levers to minimize gap-to-goal while also minimizing costs across the World-Wide Amazon workforce



Environment:

$$p(s' | s, a)_{(int+ext\ factors)}, \Delta t = 1 \text{ month}$$

State:

Talent team size, job type, job level, location, number of open jobs, market data, talent movement forecasts

Reward Function:

$$r_t = r_t^{gap} + r_t^{cost} = \beta_1 |h_{EoY} - h_{target}| - \beta_2 c_t$$

Actions:

4 actions: (jobs to open, promos, compensation, incentives)

Workforce MDP Simulator

- **Model used to define next states:** Forecast monthly talent movement based on historical trends and market indicators

Graph Transformer for Workforce Planning

Jérôme Curakus
People Experience and Technology
curakj@amazon.com

Jeanne Righy
People Experience and Technology
righyje@amazon.com

Abstract

We present a Graph Transformer deep learning method for workforce planning which can identify potential talent risks and forecast individual monthly talent movements in each segment of the Amazon corporate population. This method outperforms last-period's actual usage at Amazon by 40% in 76% of the segments and outperforms all other methods in 70% of the segments (and 17-79% of the segments). Given the dependency Graph in the proposed method can be used to interpret and understand talent forecasts, there is little drawback to using this method compared to using linear trailing rates for workforce planning.

1 Introduction

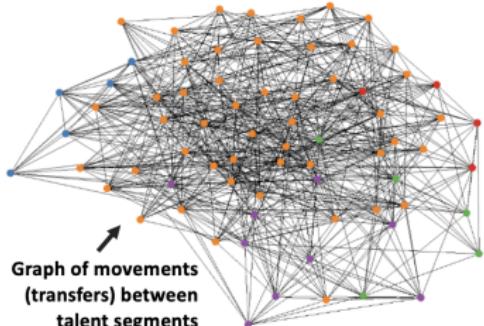
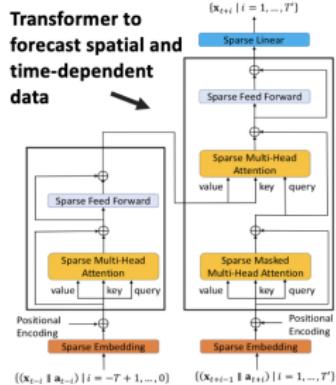
Workforce planning at Amazon sets year-end headcount targets by financial cost center to meet the company's current and future staffing needs based on business goals and talent movement forecasts (hires, promotions, transfers, attritions). Individual team leaders further plan their workforce needs by individual team, job type, job level and location. Failure to accurately forecast future headcounts and staffing needs often results in delays in productivity and costly resource allocation [1].

Forecasting Amazon talent movements is challenging due to complex spatial and temporal dependencies within the Amazon population, and non-stationarity that result from unusual events such as the Covid pandemic [2]. Amazon is an especially difficult forecasting problem due to its diverse operations worldwide and unprecedented size (over 1.6M employees in peak season [2]).

Examples of talent movement dependency at Amazon include talents in similar environments, with similar profiles and job market opportunities, or under similar talent management strategies. Many other factors, including external factors such as large swings in the company's stock value [3], can influence talent flows and thereby create correlated traffic patterns in the Amazon population.

Forecasting spatial and time-dependent data in large, complex traffic networks was recently addressed by estimating a dependency Graph that parsimoniously represents the spatial dependency between different locations in the network (nodes), and using the Graph to sparsify a deep RNN and CNN and improve the learning of long-range temporal dependency [5]. By assigning each neuron of the Transformer with a spatial location and using knowledge from the dependency Graph to prune neural connections that are not dependent, the resulting Graph Transformer could efficiently capture both spatial and temporal dependencies and scale to large traffic patterns in the Amazon population [4].

In this paper, we use a multivariate Gaussian approximation to find the dependency Graph of talent movements over the different teams of the Amazon corporate population defined by leader, job-type, and job level. We use this Graph to derive insights into the overall dynamics of Amazon talent

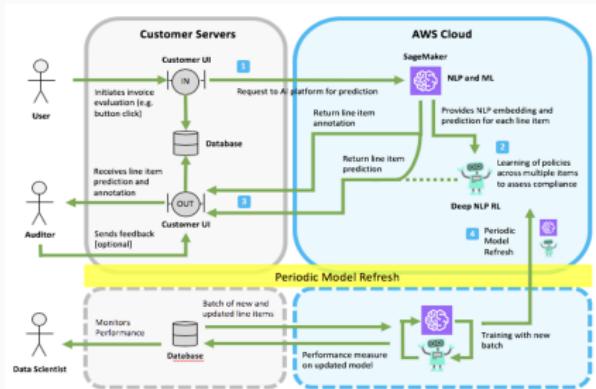


Model	MAPE	Nodes improved	MAPE in nodes improved
Graph Transformer	84%	N/A (<i>self</i>)	93%
Prophet	95%	70%	108%
s-ARIMA/State Space	87%	61%	107%
Trailing 3-month	112%	76%	133%

Audit Financial Claims with NLP

Audit Claims with Natural Language

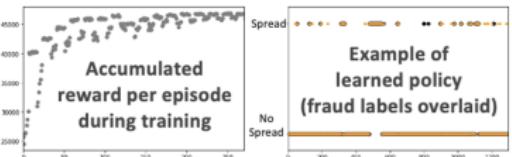
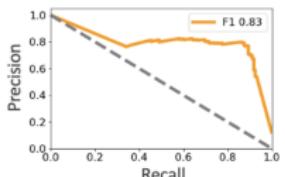
- **Goal:** Recommend compliance level of items in financial claims, to reduce time spent by human contractors, and to reduce errors



Agent finds 86% of known frauds at 80% precision

Fraud Test Results:

Precision: 80%
Recall: 86%
F1 Score: 0.83
Brier Score: 0.04
AUC Score: 0.95



Environment:

50,000 claims with at least 2 items per claim

Reward Function:

- Compliance category predicted by NLP classifier trained on individual items (regular audit)
- Frauds detected by specialized auditors

State:

- Last n items in claims encoded by NLP
- Claim metadata (source, activity code, billed units, ...)

Actions:

2 actions: Fraud risk level and compliance category

Audit Claims with Natural Language

- ▶ Post-processing to interpret RL results: Clustering in NLP space of items at risk can help auditors identify patterns of frauds more quickly

Rearrange non-compliant items per date/source/claim

	Source	Claim ID	Date	Description
Mr. Z	10113598 10113599 10113596		0000	Request and review of the end-Administrative Design of Radiopharmaceutical, Notice of Production is Required to Defendant's Interrogatories
				Request and review of all documents/replies to Plaintiff's Unverified Answers to Defendants' Interrogatories
				Request and review of the Adminstrative Design of Radiopharmaceutical, Notice of Production is Required to Plaintiff's Request for Production to Plaintiff
Miss. Y	10113600 10113600 10113600 10113600		0000	Request and review of all documents/replies to Plaintiff's Request for Production to Plaintiff
				Request and review of the Adminstrative Design of Radiopharmaceutical, Notice of Production is Required to Plaintiff's Request for Production to Plaintiff
				Request and review of all documents/replies to Plaintiff's Request for Production to Plaintiff
				Request and review of all documents/replies to Plaintiff's Request for Production to Plaintiff
Mr. X	#1731 #1732 #1733 #1734 #1735 #1736 #1737		01/15/2019 01/26/2019 10/30/2019 0000 0000 0000 0000	Request medical discovery responses in preparation of Notice of Compliance to Request for Production to Plaintiff, Requests for defendant's responses to Plaintiff's Request for Production, Requests for Plaintiff's witnesses' responses to Plaintiff's Request for Production to Plaintiff, Requests for Plaintiff's documents in preparation of Notice of Compliance to Request for Production to Plaintiff, Requests for Plaintiff's witness testimony in preparation of Notice of Compliance to Request for Production to Plaintiff, Requests for Plaintiff's documents in preparation of Notice of Compliance to Request for Production to Plaintiff, Requests for Plaintiff's witness testimony in preparation of Notice of Compliance to Request for Production to Plaintiff
				Request medical discovery responses in preparation of Notice of Compliance to Request for Production to Plaintiff, Requests for defendant's responses to Plaintiff's Request for Production, Requests for Plaintiff's witnesses' responses to Plaintiff's Request for Production to Plaintiff, Requests for Plaintiff's documents in preparation of Notice of Compliance to Request for Production to Plaintiff, Requests for Plaintiff's witness testimony in preparation of Notice of Compliance to Request for Production to Plaintiff, Requests for Plaintiff's documents in preparation of Notice of Compliance to Request for Production to Plaintiff, Requests for Plaintiff's witness testimony in preparation of Notice of Compliance to Request for Production to Plaintiff
				Request medical discovery responses in preparation of Notice of Compliance to Request for Production to Plaintiff, Requests for defendant's responses to Plaintiff's Request for Production, Requests for Plaintiff's witnesses' responses to Plaintiff's Request for Production to Plaintiff, Requests for Plaintiff's documents in preparation of Notice of Compliance to Request for Production to Plaintiff, Requests for Plaintiff's witness testimony in preparation of Notice of Compliance to Request for Production to Plaintiff, Requests for Plaintiff's documents in preparation of Notice of Compliance to Request for Production to Plaintiff, Requests for Plaintiff's witness testimony in preparation of Notice of Compliance to Request for Production to Plaintiff
				Request medical discovery responses in preparation of Notice of Compliance to Request for Production to Plaintiff, Requests for defendant's responses to Plaintiff's Request for Production, Requests for Plaintiff's witnesses' responses to Plaintiff's Request for Production to Plaintiff, Requests for Plaintiff's documents in preparation of Notice of Compliance to Request for Production to Plaintiff, Requests for Plaintiff's witness testimony in preparation of Notice of Compliance to Request for Production to Plaintiff, Requests for Plaintiff's documents in preparation of Notice of Compliance to Request for Production to Plaintiff, Requests for Plaintiff's witness testimony in preparation of Notice of Compliance to Request for Production to Plaintiff
				Request medical discovery responses in preparation of Notice of Compliance to Request for Production to Plaintiff, Requests for defendant's responses to Plaintiff's Request for Production, Requests for Plaintiff's witnesses' responses to Plaintiff's Request for Production to Plaintiff, Requests for Plaintiff's documents in preparation of Notice of Compliance to Request for Production to Plaintiff, Requests for Plaintiff's witness testimony in preparation of Notice of Compliance to Request for Production to Plaintiff, Requests for Plaintiff's documents in preparation of Notice of Compliance to Request for Production to Plaintiff, Requests for Plaintiff's witness testimony in preparation of Notice of Compliance to Request for Production to Plaintiff
				Request medical discovery responses in preparation of Notice of Compliance to Request for Production to Plaintiff, Requests for defendant's responses to Plaintiff's Request for Production, Requests for Plaintiff's witnesses' responses to Plaintiff's Request for Production to Plaintiff, Requests for Plaintiff's documents in preparation of Notice of Compliance to Request for Production to Plaintiff, Requests for Plaintiff's witness testimony in preparation of Notice of Compliance to Request for Production to Plaintiff, Requests for Plaintiff's documents in preparation of Notice of Compliance to Request for Production to Plaintiff, Requests for Plaintiff's witness testimony in preparation of Notice of Compliance to Request for Production to Plaintiff
				Request medical discovery responses in preparation of Notice of Compliance to Request for Production to Plaintiff, Requests for defendant's responses to Plaintiff's Request for Production, Requests for Plaintiff's witnesses' responses to Plaintiff's Request for Production to Plaintiff, Requests for Plaintiff's documents in preparation of Notice of Compliance to Request for Production to Plaintiff, Requests for Plaintiff's witness testimony in preparation of Notice of Compliance to Request for Production to Plaintiff, Requests for Plaintiff's documents in preparation of Notice of Compliance to Request for Production to Plaintiff, Requests for Plaintiff's witness testimony in preparation of Notice of Compliance to Request for Production to Plaintiff

...then cluster items per narrative

Example of 16 items at risk claimed by Mr. X:

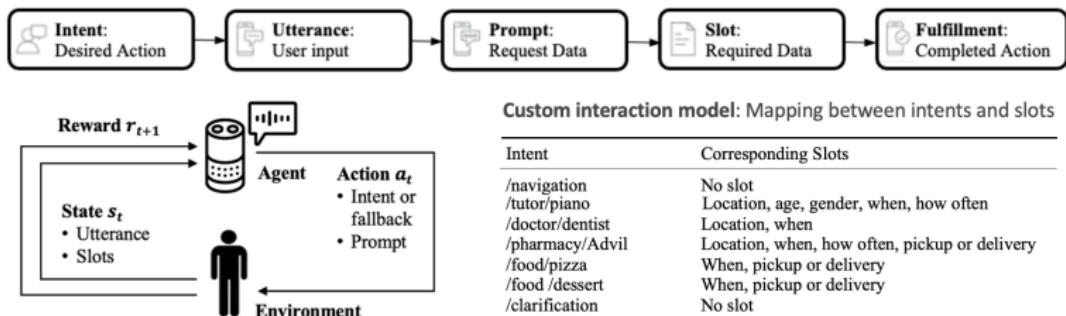
Source	Claim ID	Date	Description	Cluster ID
Mr. X	#1731	01/15/2019 0000	Initiate proper disclosure of claimant's sworn testimony (authorized by claims representative) through preparation of Notice according to S. C. Code (legal document requiring signature of S. C. attorney)	1
			sworn testimony (authorized by claims representative) on opposing counsel through preparation of statutorily required document according to S. C. Code	1
	#1732	01/26/2019 0000	Execute Amended Notice of Claimant's sworn testimony (authorized by claims representative) on opposing counsel through preparation of statutorily required document according to S. C. Code	1
			Initiate proper disclosure of claimant's sworn testimony (authorized by claims representative) through preparation of Amended Notice according to S. C. Code (legal document requiring signature of S. C. attorney)	1
		10/30/2018 0000	Advise/ment of medical discovery directive (to pursuant to - 07-214 non-applicability of HIPAA Privacy Rule to Workers' Compensation matters and applicable state law for evidentiary document production	1
			Advise/ment of medical discovery directive (to pursuant to - 07-214 non-applicability of HIPAA Privacy Rule to Workers' Compensation matters and applicable state law for evidentiary document production	1
			Advise/ment of medical discovery directive to EMS (pursuant to - 07-214 non-applicability of HIPAA Privacy Rule to Workers' Compensation matters and applicable state law for evidentiary document production	1
			Execute medical discovery request on through preparation of statutorily required document according to , and	3
			Initiate medical discovery on EMS through preparation of medical evidence request according to Reg. legal document in S. C.) to obtain pertinent documents for analysis	1
			Initiate medical discovery on through preparation of medical evidence request according to Reg. legal document in S. C. to obtain pertinent documents for analysis	2
		10/31/2018 0000	Initiate proper disclosure of claimant's sworn testimony (authorized by claims representative) through preparation of Notice according to S. C. Code (legal document requiring signature of S. C. attorney)	1
			sworn testimony (authorized by claims representative) on opposing counsel through preparation of statutorily required document according to S. C. Code	1
		13/12/2018 0000	Advise/ment of discovery directive to Claimant e . , (pursuant to - 07-214 representation of employer/carer and production of evidentiary documents (L36) (A10))	1
			Initiate discovery on claimant through preparation of subpoena evidence request according to Reg. legal document in S. C.) to obtain pertinent documents for analysis, including tax returns and business records	1
		12/05/2018 0000	Initiate proper disclosure of claimant's sworn testimony (authorized by claims representative) through preparation of Notice according to S. C. Code (legal document requiring signature of S. C. attorney)	1
			sworn testimony (authorized by claims representative) on opposing counsel through preparation of statutorily required document according to S. C. Code	1

Adaptive Dialogue Model for Chatbots

Increase Quality and Efficiency of Dialogues

- **Goal:** Identify intents, slots and fallback to maximize quality & efficiency of dialogues (Nash equilibrium in two-agent dialogues)

Components of an “interaction model” in text-speech interfaces. Example: Alexa Skills Kit



Environment:

- Library of intents with associated utterances/slots
- Single agent: Users select utterances randomly
- Two agents: Each agent becomes the environment from the perspective of the other agent

Reward Function:

- Every step: -1
- Ask for clarification: 0
- Intent and slot guesses: -5 (invalid), +5 (valid)
- Final win or loss: -10 (not happy), +10 (happy)

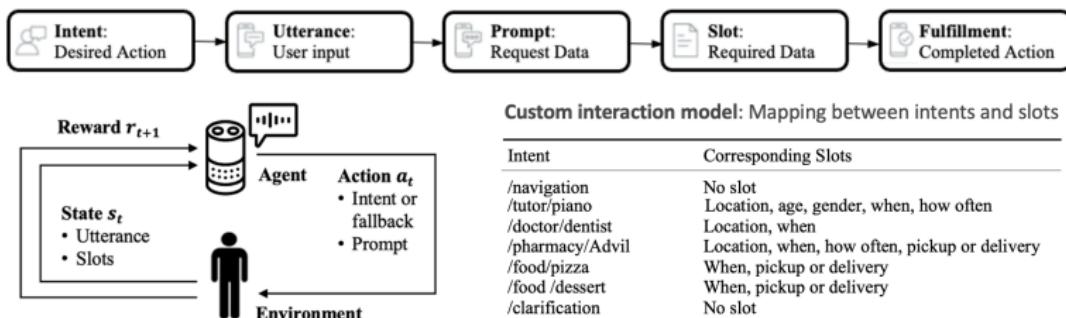
Sample of the complete custom interaction model

```
[{"intent": "/clarification",
 "utterances": ["{}"],
 "responses": ["Sorry could you please say that again?", "I apologize, can you repeat?", "I'm not sure I understood, could you say that again?"],
 "slots": []},
 {"intent": "/navigation",
 "utterances": ["Can you help me find a piano tutor?", "I need find a piano tutor", "Open the tutoring module to find a piano tutor", "Open the piano module to find a piano tutor", "I'm looking for a piano tutor", "I'm looking for a piano instructor", "Can you find a piano instructor?"],
 "responses": ["{}"],
 "slots": [{"name": "location", "type": "text"}, {"name": "age", "type": "text"}, {"name": "gender", "type": "text"}],
 "slots": [{"name": "location", "type": "text"}, {"name": "age", "type": "text"}, {"name": "gender", "type": "text"}]},
 {"intent": "/food/pizza",
 "utterances": ["Can you help me find a pizza?", "I need find a pizza", "Open the food module to find a pizza", "Open the pizza module to find a pizza", "I'm looking for a pizza", "I'm looking for a pizza to eat", "Can you find a pizza to eat?"],
 "responses": ["{}"],
 "slots": [{"name": "module", "type": "text"}, {"name": "type", "type": "text"}],
 "slots": [{"name": "module", "type": "text"}, {"name": "type", "type": "text"}]},
 {"intent": "/food/dessert",
 "utterances": ["Can you help me find some cake?", "I need find a place to buy some dessert", "I'm looking for a dessert", "I'm looking for a dessert to eat", "Can you find some cake for dessert?", "Can you find some dessert for dessert?"],
 "responses": ["{}"],
 "slots": [{"name": "module", "type": "text"}, {"name": "type", "type": "text"}],
 "slots": [{"name": "module", "type": "text"}, {"name": "type", "type": "text"}]}]
```

Increase Quality and Efficiency of Dialogues

- **Goal:** Identify intents, slots and fallback to maximize quality & efficiency of dialogues (Nash equilibrium in two-agent dialogues)

Components of an “interaction model” in text-speech interfaces. Example: Alexa Skills Kit



Environment:

- Library of intents with associated utterances/slots
- Single agent: Users select utterances randomly
- Two agents: Each agent becomes the environment from the perspective of the other agent

Reward Function:

- Every step: -1
- Ask for clarification: 0
- Intent and slot guesses: -5 (invalid), +5 (valid)
- Final win or loss: -10 (not happy), +10 (happy)

State:

- Agent 1: Utterance encoded by NLP + slots filled
- Agent 2: Prompt encoded by NLP

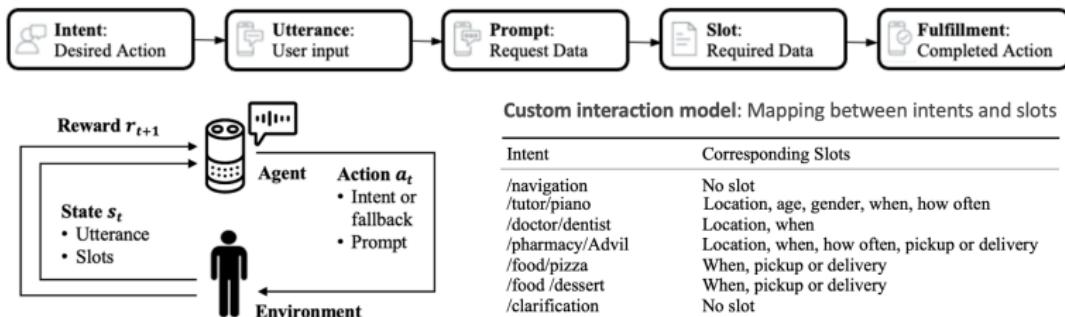
Actions:

- Agent 1: 2 actions
 - One of the 6 intents, or fallback
 - One of the 6 slots, or no slot
- Agent 2: 1 action
 - One of the utterances within intent

Increase Quality and Efficiency of Dialogues

- **Goal:** Identify intents, slots and fallback to maximize quality & efficiency of dialogues (Nash equilibrium in two-agent dialogues)

Components of an “interaction model” in text-speech interfaces. Example: Alexa Skills Kit



Environment:

- Library of intents with associated utterances/slots
- Single agent: Users select utterances randomly
- Two agents: Each agent becomes the environment from the perspective of the other agent

Reward Function:

- Every step: -1
- Ask for clarification: 0
- Intent and slot guesses: -5 (invalid), +5 (valid)
- Final win or loss: -10 (not happy), +10 (happy)

State:

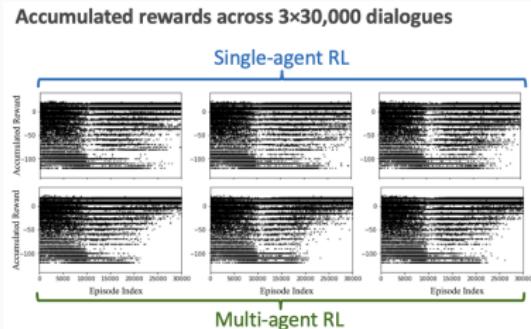
- **Agent 1:** Utterance encoded by NLP + slots filled
- **Agent 2:** Prompt encoded by NLP

Actions:

- **Agent 1:** 2 actions
 - One of the 6 intents, or fallback
 - One of the 6 slots, or no slot
- **Agent 2:** 1 action
 - One of the utterances within intent

Cooperation leads to better dialogues...

- **Analysis of Results:** Chatbot converged to policies which fulfilled intents in 99% of dialogues, in 1.8 steps on average. When users cooperated, the correct intent was fulfilled in 1.3 steps on average in 100% of dialogues



Quality and efficiency of sampled dialogues

Analysis of successful dialogues (intent filled within 10 steps)

	Single RL 0-5K	Multi RL 0-5K	Single RL 25-30K	Multi RL 25-30K
% of successful dialogues	68 (.8)	99 (.1)	67 (1.4)	100 (0)
Number of steps in successful dialogues	4.1 (.2) 4.2 (.1) 4.2 (.0)	1.0 (.0) 1.8 (.2) 1.6 (.2)	4.1 (.1) 4.1 (.1) 4.3 (.0)	1.0 (.0) 1.3 (.1) 1.3 (.0)
/navigation				
/piano				
/dentist				
/pizza				
/Advil				
/dessert				



Environment:

- Library of intents with associated utterances/slots
- Single agent: Users select utterances randomly
- Two agents: Each agent becomes the environment from the perspective of the other agent



Reward Function:

- Every step: -1
- Ask for clarification: 0
- Intent and slot guesses: -5 (invalid), +5 (valid)
- Final win or loss: -10 (not happy), +10 (happy)



State:

- **Agent 1:** Utterance encoded by NLP + slots filled
- **Agent 2:** Prompt encoded by NLP



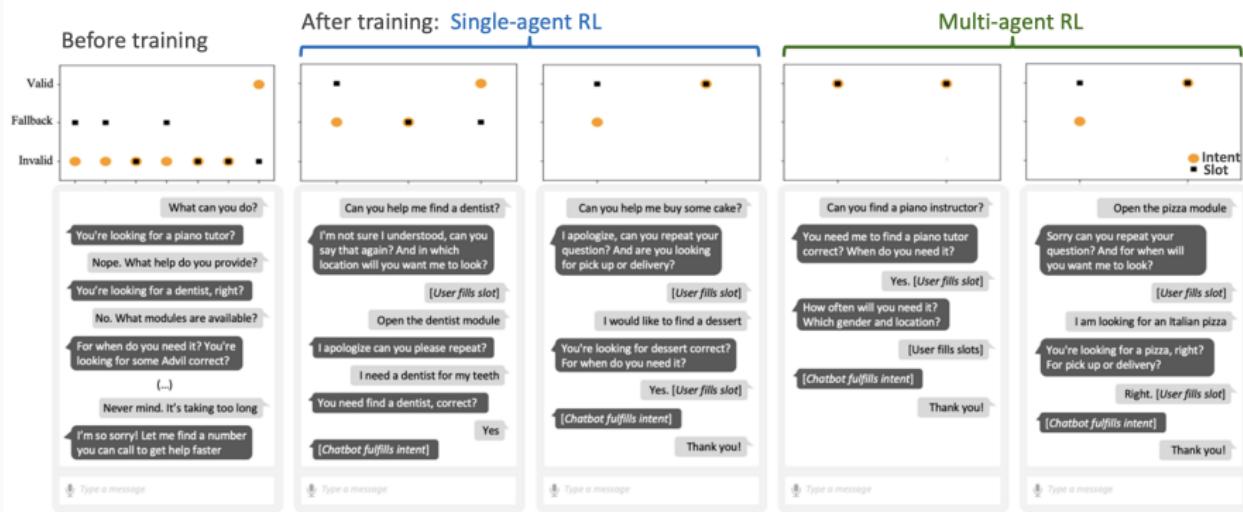
Actions:

- **Agent 1:** 2 actions
 - One of the 6 intents, or fallback
 - One of the 6 slots, or no slot
- **Agent 2:** 1 action
 - One of the utterances within intent

The chatbot learned an original strategy...

- **Analysis of Results:** The chatbot identified original strategies to increase speed of fulfillment without sacrificing coherence, such as filling in valid slots even when utterances are too ambiguous to identify the exact intent

Sample of user-bot interactions. Validity of actions taken by the chatbot at every step is shown in the top caption



Your Turn!