Motivation	Simulation	Algorithms	Experiments	Lessons Learned	Next Steps
000000	00000000000	0000000000000	00000000	00000000	000

# Communicating Agents Seeking Information

David James

2019-01-28

David James Communicating Agents Seeking Information

Motivation	Algorithms	Experiments	Lessons Learned	Next Steps













Motivation	Simulation	Algorithms	Experiments	Lessons Learned	Next Steps
●00000	00000000000	0000000000000	00000000	00000000	

#### Table of Contents

- 1 Motivation
- 2 Simulation
- 3 Algorithms
- 4 Experiments
- 5 Lessons Learned
- 6 Next Steps

#### Personal Motivations

I am fascinated by:

- simulation as a tool
- building agents that learn
- understanding multi-agent systems
- $\bullet$  understanding how communication affects  $\uparrow$

My educational background:

- Liberal Arts + ECE (undergraduate)
- Public Policy + CS (graduate)

Motivation 00●000	Simulation 000000000000	Algorithms 0000000000000	Experiments 00000000	Lessons Learned 00000000	Next Steps 000
_					

Research Question

How do decentralized\* communicating agents behave?

- What behavior patterns emerge?
- How do incentives affect behavior?
- Do agents reach global optima?
- What are the algorithmic challenges?

Motivation	Simulation	Algorithms	Experiments	Lessons Learned	Next Steps
000●00	000000000000	0000000000000	00000000		000
Design	Process				

What kind of simulation can explore these questions?

Let's walk through some design questions:

- How does communication work?
- How does learning (optimization) work?

Motivation	Simulation	Algorithms	Experiments	Lessons Learned	Next Steps
000000					

# Communication Type



Motivation	Simulation	Algorithms	Experiments	Lessons Learned	Next Steps
00000●	00000000000	000000000000	00000000	000000000	000

#### Learning Type



#### Decentralized Learning



#### Table of Contents





#### 3 Algorithms

4 Experiments

5 Lessons Learned

#### 6 Next Steps

David James Communicating Agents Seeking Information



You can think of the CASI simulation as a game:

- There are  $n_a$  agents.
- Each episode consists of *T* steps.
- If an agent correctly guesses an unsolved goal attribute, it gains points.
- If an agent guesses all of its  $n_g$  goal attributes, it gains many points and wins.
- Various actions are penalized (more on this later).
- For example, asking a question has a small cost.

Motivation 000000	Simulation 0000000000	Algorithms 0000000000000	Experiments 00000000	Lessons Learned 00000000	Next Steps 000	
Simulation Attributes						

Each simulation has a configurable number of attributes.

Here is an example with three agents and six attributes:



Attribute Vectors

Motivation 000000	Simulation 000000000000	Algorithms 0000000000000	Experiments 00000000	Lessons Learned 00000000	Next Steps 000
Agent'	s Perspectiv	ve			

Each agent sees an observation and must choose an action.

An observation consists of:

- the agent's internal state
- interactions with other agents
  - (answers received, questions asked)

An action consists of:

- intentions on how to respond to questions
- choosing what to do next, one of:
  - (do nothing, guess an attribute, ask a question)

#### Communication Protocol

- agent *i* says to agent *j*:
   "tell me about attribute *a*"
- agent j chooses an intention (i.e. to lie, ignore, or tell the truth)
- the simulation converts the intention to an answer and delivers it to agent *i*

Motivation	Simulation	Algorithms	Experiments	Lessons Learned	Next Steps
000000	00000000000	0000000000000	00000000	00000000	000

#### Agent Intentions

Intention	Response		
	if <i>unknown</i>	if <i>known</i>	
	i.e. <i>None</i>	i.e. <i>Some</i> (v)	
specific lie	Value(r)	$Value(r) \mid r \neq v$	
general lie	Known	Unknown	
ignore	None	None	
incomplete truth	Unknown	Known	
complete truth	Unknown	Value(v)	

where r is a random value

Motivation 000000	Simulation 000000000000	Algorithms 0000000000000	Experiments 00000000	Lessons Learned 00000000	Next Steps 000
Incent	ives				

```
let incentives = Incentives { // Experiment #1
    cost_ask_question: 2.0,
    cost_ask_self_question: 20.0,
    cost_excessive_guessing: 20.0,
    cost_incorrect_guess: 20.0,
    cost_known_guess: 20.0,
    cost_non_goal_guess: 20.0,
    cost_solved_guess: 20.0,
    cost_unnecessary_reply: 20.0,
    max_guesses_per_attr: 1,
    reward_correct: 100.0,
    reward_win: 200.0,
}:
```

David James Communicating Agents Seeking Information

Motivation	Simulation	Algorithms	Experiments	Lessons Learned	Next Steps
000000	0000000●0000	0000000000000	00000000	00000000	000

#### CASI Architecture



- Simulator (written in Rust)
- Learning Agents (written with PyTorch)
- Interprocess communication (uses gRPC)

Motivation 000000	Simulation 0000000000000	Algorithms 0000000000000	Experiments 00000000	Lessons Learned 00000000	Next Steps 000
Why R	ust?				



Rust helps me write (and refactor) fast and correct software, quickly.

- Performance: no runtime, no garbage collector
- *Reliability*: expressive type system gives memory-safety, thread-safety
- Productivity: package manager, libraries, community



# Ċ

As you might expect, PyTorch:

- manages a computation graph and auto-differentiation
- includes reusable modules for neural networks
- has CPU and GPU support

In comparison to TensorFlow, I find PyTorch is:

- more natural (has a more intuitive API)
- easier to debug



# **GRPC**

Simple service definition:

• easy refactoring

Works across languages and platforms:

- conveniently bridges Rust and Python
- public and private use cases (RPC, APIs)

Start quickly and scale:

- works on one machine
- viable for large multi-agent systems

 Motivation
 Simulation
 Algorithms
 Experiments
 Lessons
 Learned
 Next
 Steps

 000000
 0000000000
 000000000
 00000000
 000
 000
 000
 000
 000
 000
 000
 000
 000
 000
 000
 000
 000
 000
 000
 000
 000
 000
 000
 000
 000
 000
 000
 000
 000
 000
 000
 000
 000
 000
 000
 000
 000
 000
 000
 000
 000
 000
 000
 000
 000
 000
 000
 000
 000
 000
 000
 000
 000
 000
 000
 000
 000
 000
 000
 000
 000
 000
 000
 000
 000
 000
 000
 000
 000
 000
 000
 000
 000
 000
 000
 000
 000
 000
 000
 000
 000
 000
 000
 000
 000

#### Neural Network gRPC Definition

Here is how I use gRPC to define a neural network model for a learning agent:

```
// Neural network for action prefs and value function
service Net {
    rpc Init (NetCfg) returns (Empty) {}
    rpc InitHidden (Empty) returns (Empty) {}
    rpc Predict (Observations) returns (Predictions) {}
    rpc ResetOptimizer (OptimCfg) returns (Empty) {}
    rpc Train (Sequences) returns (Loss) {}
    rpc GetParams (Empty) returns (Params) {}
    rpc SetParams (Params) returns (Empty) {}
}
```

This file will generate executable code in a wide variety of languages.

Motivation	Simulation	Algorithms	Experiments	Lessons Learned	Next Steps
000000	00000000000	•00000000000	00000000	000000000	000

#### Table of Contents

- Motivation
- 2 Simulation
- 3 Algorithms
- 4 Experiments
- 5 Lessons Learned
- 6 Next Steps



Reinforcement learning (RL) is learning what to do how to map situations to actions — so as to maximize a numerical reward signal. - Sutton & Barto



Motivation 000000	Simulation 00000000000	Algorithms 000000000000	Experiments 00000000	Lessons Learned 000000000	Next Steps 000
Decent	ralized RL				

You may remember this diagram from before:



Can you think about some pros and cons?

Motivation 000000	Simulation 000000000000	Algorithms 000●000000000	Experiments 00000000	Lessons Learned 00000000	Next Steps 000
	mensions				

Environmental considerations:

- states: fully-observed vs. partially-observed
- actions: discrete vs. continuous

Algorithmic considerations:

- policies: deterministic vs. stochastic
- on-policy vs. off-policy
- model-based vs. model-free \*

Motivation 000000	Simulation 00000000000	Algorithms 0000●00000000	Experiments 00000000	Lessons Learned 00000000	Next Steps
'Model	s' in RL				

I briefly wanted to mention what 'model' means in the context of RL, since it can be confusing.

Ben Recht says it well in "A Tour of Reinforcement Learning: The View from Continuous Control" (2018):

The term "model-free" almost always means "no model of the state transition function" when casually claimed in reinforcement learning research. However, this does not mean that modeling is not heavily built into the assumptions of model-free RL algorithms.

Motivation	Simulation	Algorithms	Experiments	Lessons Learned	Next Steps
000000	00000000000	000000000000	00000000	00000000	000

#### Algorithm Experimentation

I've implemented these algorithms from scratch in Rust and PyTorch:

- REINFORCE ('vanilla')
- REINFORCE with baseline
- Proximal Policy Optimization (PPO)
- Advantage Actor Critic (A2C)

I'm currently using A2C + PPO with exploration bonuses.

Motivation	Algorithms	Experiments	Lessons Learned	Next Steps
	00000000000			

# REINFORCE

#### **REINFORCE:** Monte-Carlo Policy-Gradient Control (episodic) for $\pi_*$

 $\begin{array}{l} \mbox{Input: a differentiable policy parameterization } \pi(a|s, \pmb{\theta}) \\ \mbox{Algorithm parameter: step size } \alpha > 0 \\ \mbox{Initialize policy parameter } \pmb{\theta} \in \mathbb{R}^{d'} (\mbox{e.g., to } \mathbf{0}) \\ \mbox{Loop forever (for each episode):} \\ \mbox{Generate an episode } S_0, A_0, R_1, \dots, S_{T-1}, A_{T-1}, R_T, \mbox{ following } \pi(\cdot|\cdot, \pmb{\theta}) \\ \mbox{Loop for each step of the episode } t = 0, 1, \dots, T-1: \\ \mbox{} G \leftarrow \sum_{k=t+1}^{T} \gamma^{k-t-1} R_k \\ \mbox{} \theta \leftarrow \pmb{\theta} + \alpha \gamma^t G \nabla \ln \pi(A_t|S_t, \pmb{\theta}) \end{array}$ 

"Reinforcement Learning: An Introduction" by Sutton & Barto, 2018

Motivation	Simulation	Algorithms	Experiments	Lessons Learned	Next Steps
000000	00000000000	0000000000000	00000000	000000000	000

#### **PPO** Actor-Critic

#### Algorithm 1 PPO, Actor-Critic Style

for iteration=1,2,... do for actor=1,2,..., N do Run policy  $\pi_{\theta_{old}}$  in environment for T timesteps Compute advantage estimates  $\hat{A}_1, \ldots, \hat{A}_T$ end for Optimize surrogate L wrt  $\theta$ , with K epochs and minibatch size  $M \leq NT$   $\theta_{old} \leftarrow \theta$ end for

"Proximal Policy Optimization Algorithms" by Schulman et al. 2017

Motivation 000000	Simulation 000000000000	Algorithms 0000000●0000	Experiments 00000000	Lessons Learned 00000000	Next Steps 000
Explor	ration				

There are many kinds of exploration in RL:

- $\epsilon$ -greedy exploration \*
- a result of stochastic policies
- entropy-based exploration
- seeking novel (state, action) pairs

# Count-Based Exploration with Hashing

From #Exploration by Tang et. al, 2017:

a simple generalization of the classic count-based approach [...]. States are mapped to hash codes, which allows to count their occurrences with a hash table. These counts are then used to compute a reward bonus according to the classic count-based exploration theory.

# Locality Sensitive Hashing (LSH)

From Wikipedia:

LSH hashes input items so that similar items map to the same 'buckets' with high probability [...]. LSH differs from conventional and cryptographic hash functions because it aims to maximize the probability of a 'collision' for similar items.

Arguably, the key criteria for choosing an LSH algorithm is the relevant *distance metric*.

#### LSH for Hamming Distance

Early LSH papers on Hamming distances suggest using single-column sampling techniques.

- Each output bit is created simply by sampling a single column from the input vector.
- However, sampling only a fraction of the input vector seems suboptimal. I would expect benefits from drawing information from all columns.

#### Hamming Distance LSH with Super-Bits

- I blended two approaches from the literature as follows:
  - Sample from each column but only once. Otherwise, differences can be accentuated.
  - Designate groups of the output bits. This allows trading-off between the number of 'anchor' points for the distance calculation vs. the amount of compression. (Inspired by Super-Bit LSH by Jie et al., 2012.)
  - Combine columns using the Hamming distance. Store the result in each super-bit group.



#### Table of Contents

- Motivation
- 2 Simulation
- 3 Algorithms
- 4 Experiments
- 5 Lessons Learned

#### 6 Next Steps

Motivation 000000	Simulation 000000000000	Algorithms 0000000000000	Experiments 0●000000	Lessons Learned 00000000	Next Steps 000
Attribu	ites				

Here is an example of a Rust data structure to configure attributes for an experiment.

```
/// Attribute configuration
pub struct AttrCfg {
    /// Total number of attributes in simulation
    pub count: u16,
```

/// Attributes known from direct experience (per agent)
pub direct: u16,

/// To win, an agent must solve this many attributes
pub goal: u16,

/// Number of possible values for each attribute pub values: u8,

}

#### Experiment 0 Description

#### Two Agents:

- One non-learning random agent
- One learning agent
  - a feedforward neural network

Question: Can this agent learn to win?

Answer: Perhaps surprisingly, it can, by exploiting the lack of random initialization.

You might call it 'Groundhog Day' learning.



Simulation 000000000000 Algorithms

Experiments

Lessons Learner

Next Steps

#### Experiment 0 Output 1

Running experiment el		
<pre>iteration, ret_mean, ret_std, ret_</pre>	_max, ret_min, loss_mean	
201	011.000.000 1000000.1000000:000.000	-0.4073
14	011.000.001 1000000.1000000:000.000	-0.4081
70	011.000.010 1000000.1000000:000.000	-0.4056
43	011.000.011 1000000.1000000:000.000	-0.4045
31	011.000.100 1000000.1000000:000.000	-0.4105
15	011.000.101 1000000.1000000:000.000	-0.4115
16	011.000.110 1000000.1000000:000.000	-0.4048
45	011.000.111 1000000.1000000:000.000	-0.4050
71	011.001.000 1000000.1000000:000.000	-0.4068
10	011.001.010 1000000.1000000:000.000	-0.3940
9	011.001.110 1000000.1000000:000.000	-0.3880
77	011.010.000 1000000.1000000:000.000	-0.4115
26	011.010.001 1000000.1000000:000.000	-0.4117
14	011.010.101 1000000.1000000:000.000	-0.4149
121	011.011.000 1000000.1000000:000.000	-0.3929
17	011.011.100 1000000.1000000:000.000	-0.3963
[best return= -337.672]		
0, -337.672, 126.714,	3.619, -474.629, 1742.59	

Simulation

Algorithms

Experiments

Lessons Learned

Next Steps

#### Experiment 0 Output 2

3	-67.7827
4   011.000.101   1000000.1000000:000.000	-71.8345
	-73.3816
60	-70.7561
7	-79.1143
21     011.001.100 1000000.1000000:000.000	-83.2526
91    011.001.110 1000000.1000000:000.000	-81.9762
9     011.010.000 1000000.1000000:000.000	-78.2183
1    011.010.001 1000000.1000000:000.000	-75.4024
	-81.0234
	-78.2075
	-86.8079
388     011.011.100 1000000.1000000:000.000	-89.6306
[best return= 64.967]	
43, 64.967, 90.441, 217.781, -133.422, 197.47	

Motivation 000000 Simulation

Algorithms 000000000000000 Experiments

Lessons Learned

Next Steps

#### Experiment 0 Output 3



#### Experiment 1 Description

Same as Experiment 0 but with random experiment initialization:

Two Agents:

- One non-learning random agent
- One learning agent
  - a feedforward neural network

Question: How can this agent learn to win?

Answer: It can't; communication would be required.

#### Experiment 2 Description

#### Two Agents:

- One non-learning completely honest agent
- One learning agent, with both:
  - a feedforward neural network
  - an LSTM

Question: How can this agent learn to win?

Answer: By asking questions of the honest agent.

Discuss: Is the LSTM necessary in this case?

Motivation	Simulation	Algorithms	Experiments	Lessons Learned	Next Steps
000000	00000000000	0000000000000	00000000	●00000000	000

#### Table of Contents

- Motivation
- 2 Simulation
- 3 Algorithms
- 4 Experiments
- **5** Lessons Learned

#### 6 Next Steps

#### Problem: Diverse Purposes for Runs

Since the simulation engine is a platform for experimentation, it has different purposes:

- Verifying new algorithms
- Inspecting policy changes over time
- Exploring hyperparameters
- Monitoring optimizer performance

Getting the verbosity level is tricky: **too little** misses out, but **too much** creates clutter...

#### Solution: Granular Logging

One solution is to use fine-grained logging:

```
let log_cfg = LogCfg { // subset
    log_action_probs: false,
    log_actions_rewards: false,
    log_attrs: false,
    log_encodings: false,
    log_experiment_start: true,
    log_interactions: false,
    log_iteration_csv: true,
    log_learning: false,
    log_obs_summary: true,
    log_trajectories: false,
};
```

# Problem: When to Stop Optimization?

Policy gradient (PG) methods aim to improve a policy over many iterations.

However, the policy performance does not necessarily monotonically improve.

Empirically, PG methods, even with improvements such as PPO, can sometimes drift away from good performance and never recover.

See: "Where Did My Optimum Go?: An Empirical Analysis of Gradient Descent Optimization in Policy Gradient Methods" by Henderson et al. 2018.

### Why do Policy Gradient Methods Degrade?

The Policy Gradient Theorem gives the correct gradient, so what could go wrong in practice?

- Gradient ascent takes many iterations. How many? What does research say about the bounds as applied to PG?
- The underlying function approximator\* may lack sufficient representational power.
- The network parameters may not be able to efficiently respond to the latest distributional shift.

# Solution: Backtracking Algorithm

I implemented a three-level backtracking algorithm that requires three hyper-parameters:

```
/// Backtracking configuration
pub struct BacktrackCfg {
    /// Snapshot pops before stopping early
    pub pop_limit: u32,
    /// Consecutive resets before snapshot pop
    pub restart_limit: u32,
    /// Consecutive unimproved iterations before reset
    pub unimproved_limit: u32,
}
```

The algorithm is shown on the next slide.

# Backtracking Algorithm

```
Algorithm 1: Three-level backtracking algorithm
snapshots \leftarrow []; pops \leftarrow 0; restarts \leftarrow 0; unimproved \leftarrow 0
for iteration \leftarrow 0 to max iterations do
    if improved_performance then
         snapshots.push(model); restarts \leftarrow 0; unimproved \leftarrow 0
    else unimproved \leftarrow unimproved + 1
    if unimproved > max_unimproved then
         restarts \leftarrow restarts + 1; unimproved \leftarrow 0
         if restarts > max restarts then
              pops \leftarrow pops + 1
             if pops > max_pops then stop_early()
             else snapshots.pop(); restarts \leftarrow 0
              model \leftarrow snapshots.last()
    else
         model.optimization_step()
```

#### Problem: Consistency of Responses

An agent reports an intention to the simulation; the simulation converts it to a response that another agent observes.



source: http://factmyth.com

Claim: inconsistent lies are easier to discover.

How can the simulation achieve consistent lying?

Motivation 000000	Simulation 000000000000	Algorithms 0000000000000	Experiments 00000000	Lessons Learned 0000000●	Next Steps 000
Consis	stent Lving				

```
/// Return a deterministic lie that depends only on:
/// 1. the function arguments
/// 2. the agent's private salt
/// Note: 'q' / 'a' = agent id of questioner / answerer
fn specific_lie(
    &self, q: AgentId, attr_idx: u16, attr_val: Option<u8>
) -> Answer \{
   let m = self.attr_cfg.values as u64;
    let hash = self.calc_hash(q, attr_idx);
   let lie = match attr val {
        Some(true val) =>
          (true_val as u64 + (hash \% (m - 1))) \% m.
        None => hash % m
    }:
    Answer::Value(lie as u8)
}
```

Motivation	Simulation	Algorithms	Experiments	Lessons Learned	Next Steps
000000	00000000000	0000000000000	00000000	000000000	●00

#### Table of Contents

- Motivation
- 2 Simulation
- 3 Algorithms
- 4 Experiments
- 5 Lessons Learned



Motivation 000000	Simulation 000000000000	Algorithms 0000000000000	Experiments 00000000	Lessons Learned 00000000	Next Steps ○●○
Next S	teps				

- LSTM architecture & parameter tuning
- Create experiments with multiple learning agents
- Improve visualizations
- Try various multi-agent reward structures
- Use unsupervised learning to find patterns in policies
- Improve efficiency by using parallelism in Rust

Motivation 000000	Simulation 000000000000	Algorithms 0000000000000	Experiments 00000000	Lessons Learned 000000000	Next Steps ○○●
Discus	sion				

I look forward to your questions and comments.