# Learning to Win by Reading Manuals in a Monte-Carlo Framework
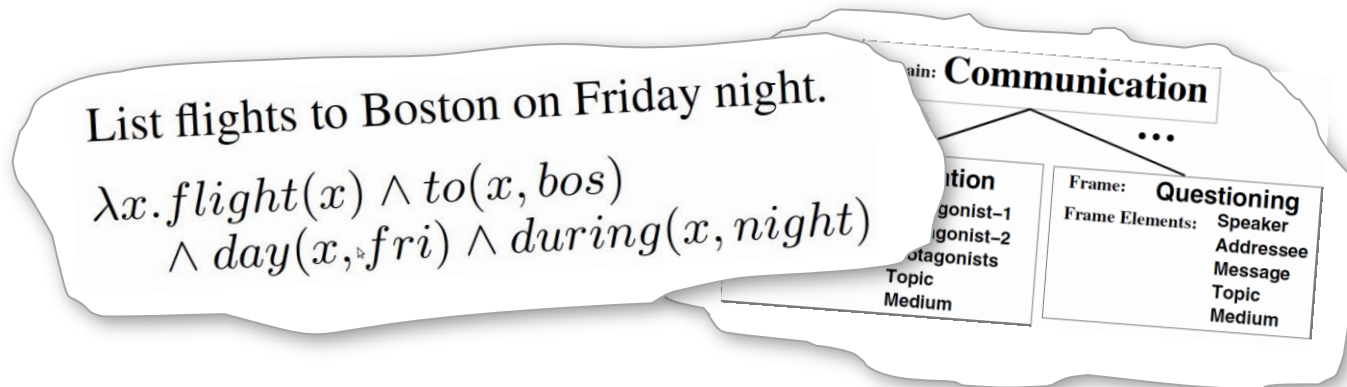
S.R.K. Branavan,  David Silver,  Regina Barzilay

MIT

# Semantic Interpretation

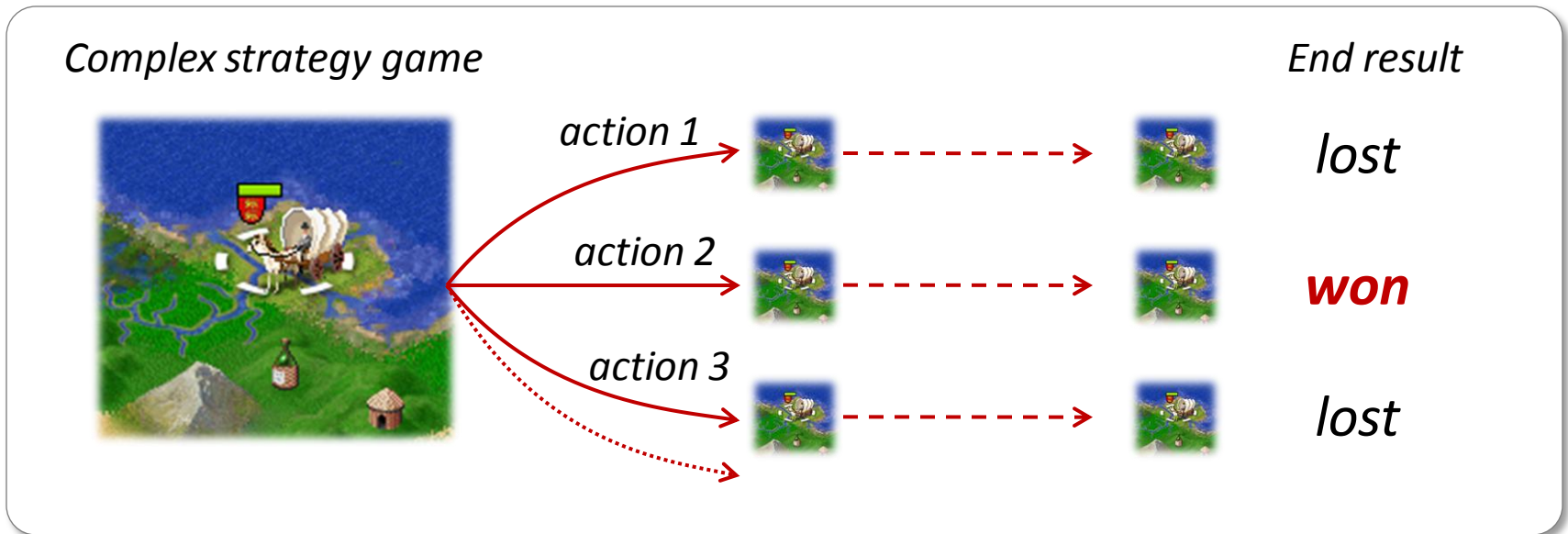Traditional view:

*Map text into an abstract representation*

List flights to Boston on Friday night.

$$\lambda x. flight(x) \wedge to(x, bos)$$
$$\wedge \, day(x, fri) \wedge during(x, night)$$

Communication

Questioning
Frame Elements: Speaker
Addressee
Message
Topic
Medium

Alternative view:

*Map text into a representation which helps performance in a control application*

# Semantic Interpretation for Control Applications



Traditional approach:

*Learn action-selection policy from game feedback.*

Our contribution:

*Use textual advice to guide action-selection policy.*

# Leveraging Textual Advice: Challenges

1. Find sentences relevant to given game state.

*Game state*



*city*    *settler*

*Strategy document*

*You start with two settler units. Although settlers are capable of performing a variety of useful tasks, your first task is to move the settlers to a site that is suitable for the construction of your first city. Use settlers to build the city on grassland with a river running through it if possible. You can also use settlers to irrigate land near your city.   In order to survive and grow …*

# Leveraging Textual Advice: Challenges

1. Find sentences relevant to given game state.

*Game state*

*Strategy document*

*city*   *settler*

*You start with two settler units. Although settlers are capable of performing a variety of useful tasks, your first task is to move the settlers to a site that is suitable for the construction of your first city.*
*Use settlers to build the city on grassland with a river running through it if possible.*
*You can also use settlers to **irrigate land near your cit**y.   In order to survive and grow ...*

# Leveraging Textual Advice: Challenges

1. Find sentences relevant to given game state.

*Game state*

*Strategy document*

city

settler

settler

*You start with two settler units. Although settlers are capable of performing a variety of useful tasks, your first task is to move the settlers to a site that is suitable for the construction of your first city.*
*Use settlers to **build the city on grassland with a river** running through it if possible.*
*You can also use settlers to irrigate land near your city.* *In order to survive and grow ...*

# Leveraging Textual Advice: Challenges

2. Label sentences with predicate stucture.

*Move* the **settler** to a site suitable for **building** a **city**, onto grassland with a river if possible.

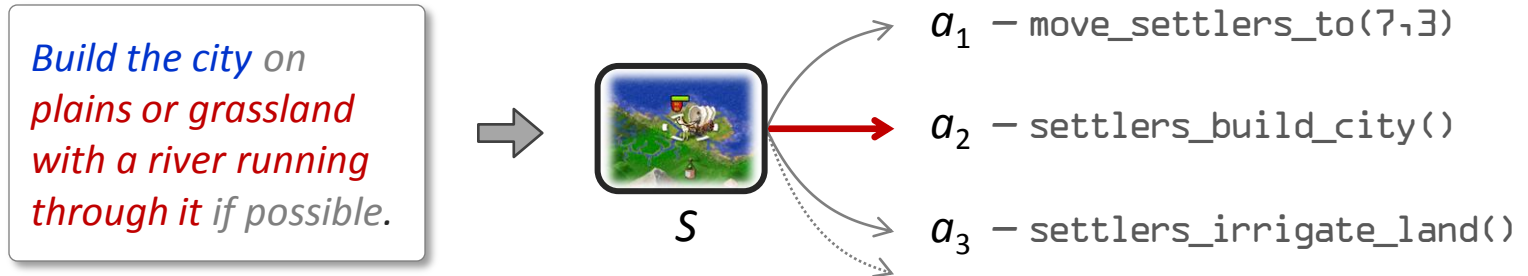`move_settlers_to()` **?**

`settlers_build_city()` **?**

*Move the settler* to a site suitable for building a city, onto grassland with a river if possible.

`move_settlers_to()`

Label words as *action*, *state* or *background*

# Leveraging Textual Advice: Challenges
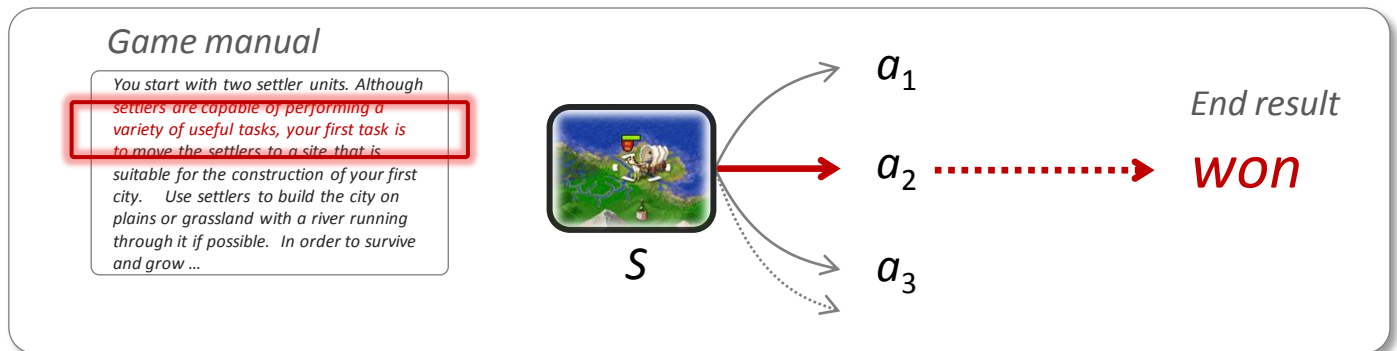
3. Guide action selection using relevant text



*Build the city* on *plains or grassland with a river running through it* if possible.

$S$

$a_1$ — move_settlers_to(7,3)

$a_2$ — settlers_build_city()

$a_3$ — settlers_irrigate_land()
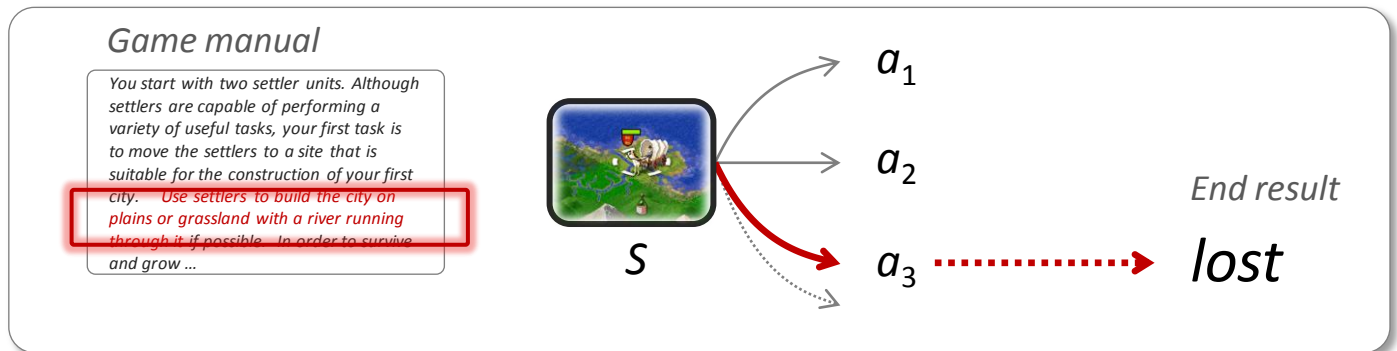
# Learning from Game Feedback

**Goal:** *Learn from game feedback as only source of supervision.*

**Key idea:** *Better parameter settings will lead to more victories.*

**Model params:** $\theta_1$



*Game manual*

You start with two settler units. Although settlers are capable of performing a variety of useful tasks, your first task is to move the settlers to a site that is suitable for the construction of your first city. Use settlers to build the city on plains or grassland with a river running through it if possible. In order to survive and grow …

$a_1$

$s$ → $a_2$ ⟶ *End result* **won**

$a_3$

**Model params:** $\theta_2$

*Game manual*

You start with two settler units. Although settlers are capable of performing a variety of useful tasks, your first task is to move the settlers to a site that is suitable for the construction of your first city. Use settlers to build the city on plains or grassland with a river running through it if possible. In order to survive and grow …

$a_1$

$a_2$

$s$ → $a_3$ ⟶ *End result* **lost**

# Model Overview

➡️ <span style="color:red">**Monte-Carlo Search Framework**</span>
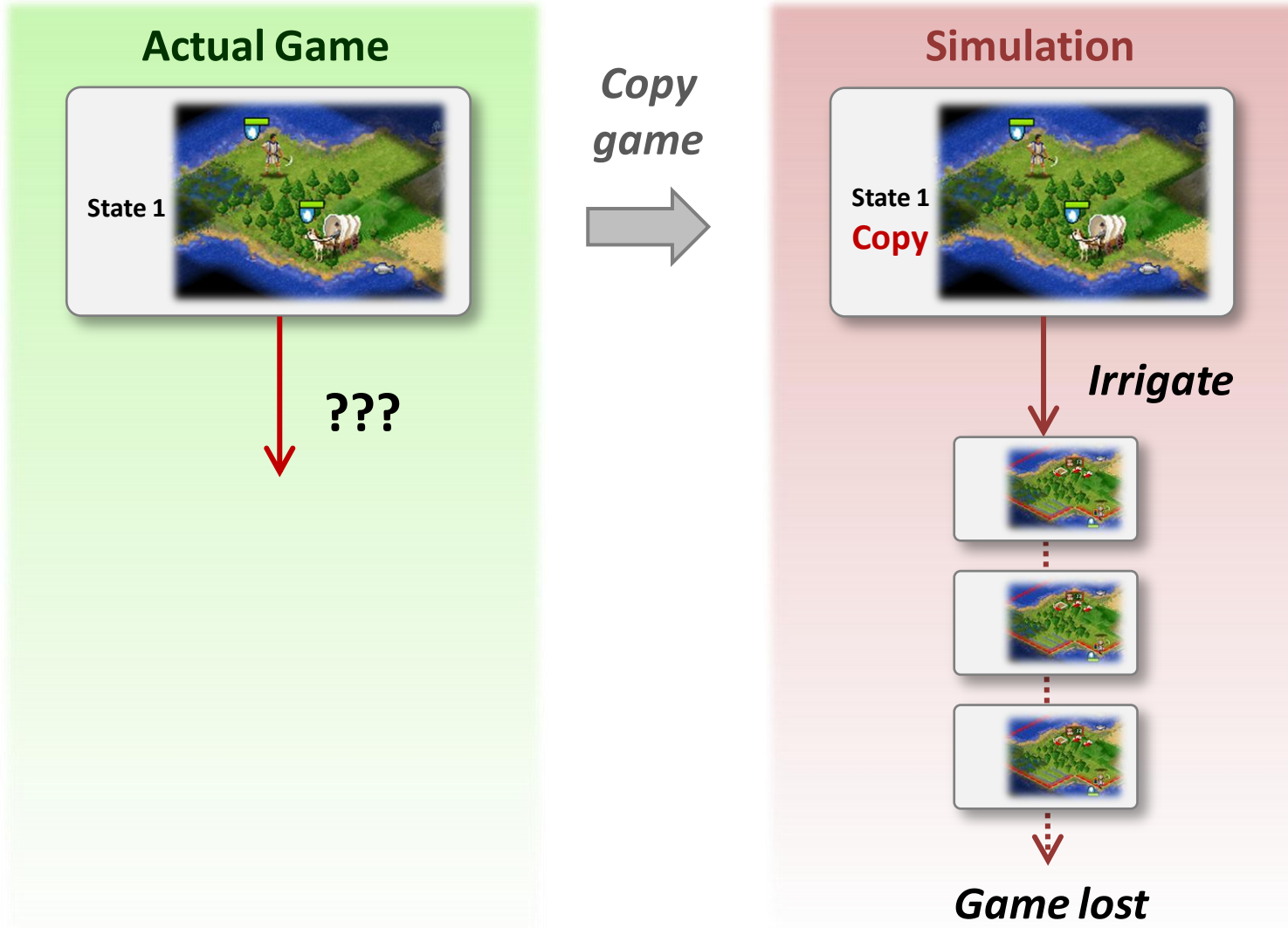
- *Learn action selection policy from simulations*

- *Very successful in complex games like Go and Poker.*

<span style="color:red">**Our Algorithm**</span>

- *Learn text interpretation from simulation feedback*
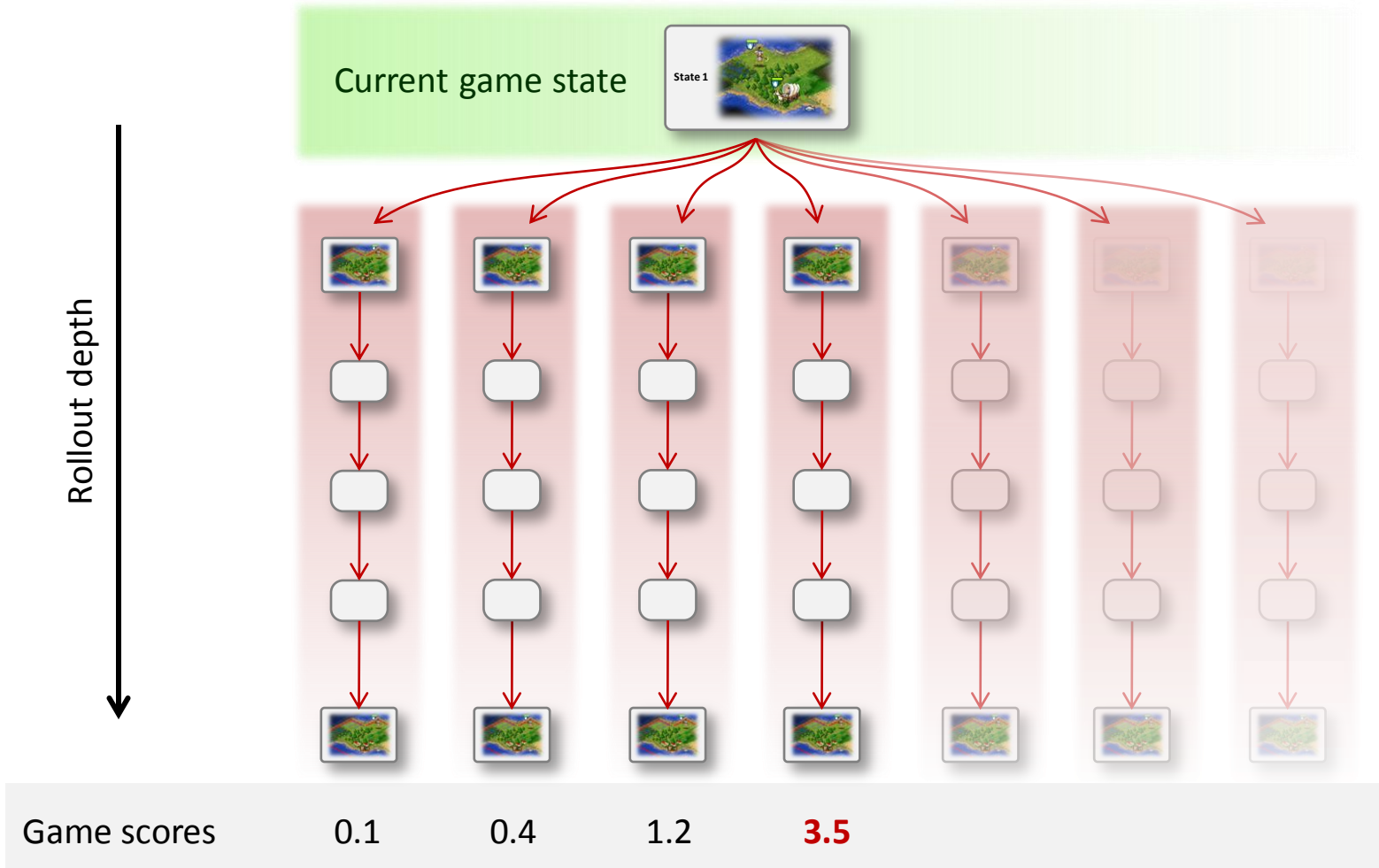
- *Bias action selection policy using text*

# Monte-Carlo Search

*Select actions via simulations, game and opponent can be stochastic*
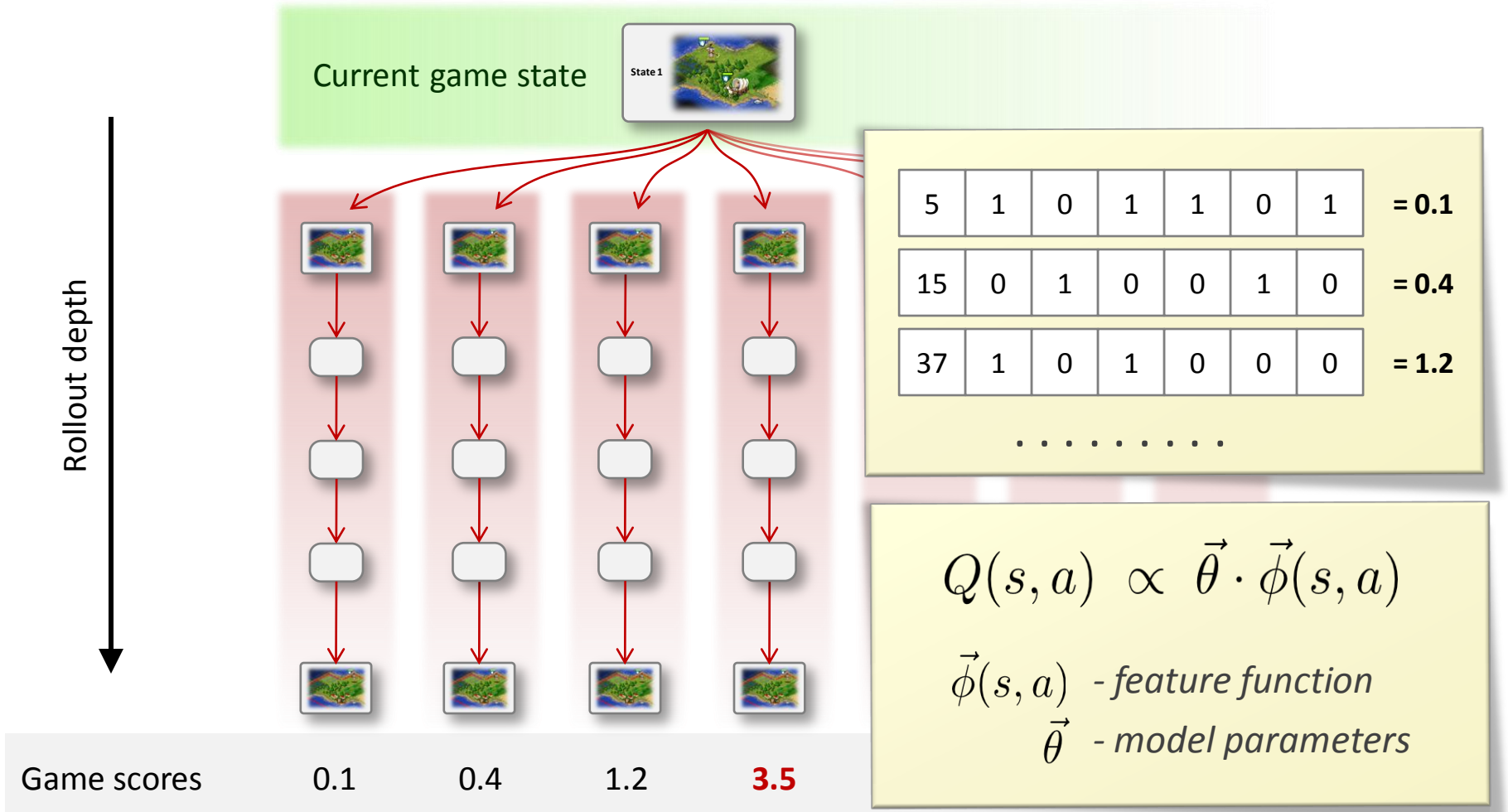
# Monte-Carlo Search

Try many candidate actions from current state & see how well they perform.



Current game state

State 1

Rollout depth

| Game scores | 0.1 | 0.4 | 1.2 | **3.5** |
|---|---|---|---|---|

# Monte-Carlo Search

Try many candidate actions from current state & see how well they perform.

Learn feature weights from simulation outcomes



Current game state

State 1

Rollout depth

| 5 | 1 | 0 | 1 | 1 | 0 | 1 | = 0.1 |

| 15 | 0 | 1 | 0 | 0 | 1 | 0 | = 0.4 |

| 37 | 1 | 0 | 1 | 0 | 0 | 0 | = 1.2 |

. . . . . . . . . .

$$Q(s,a) \propto \vec{\theta} \cdot \vec{\phi}(s,a)$$

$\vec{\phi}(s,a)$ - *feature function*

$\vec{\theta}$ - *model parameters*

Game scores    0.1    0.4    1.2    **3.5**

# Model Overview

Monte-Carlo Search Framework

- *Learn action selection policy from simulations*

Our Algorithm

- *Bias action selection policy using text*

- *Learn text interpretation from simulation feedback*

# Modeling Requirements

- *Identify sentence relevant to game state*


Build cities near rivers or ocean.

- *Label sentence with predicate structure*

Build cities near rivers or ocean. → **Build cities** *near* **rivers or ocean.**

- *Estimate value of candidate actions*


**Build cities** *near* **rivers or ocean.** →

| | |
|---|---|
| *Irrigate :* | -10 |
| *Fortify :* | -5 |
| . . . . | |
| *Build city :* | 25 |

# Sentence Relevance

*Identify sentence relevant to game state and action*

State $s$, candidate action $a$, document $d$

$$p(y = y_i | s, a, d) \propto e^{\vec{u} \cdot \vec{\phi}(y_i, s, a, d)}$$

Sentence $y_i$ is selected as relevant

Log-linear model:
$$\vec{u} \quad \text{- weight vector}$$
$$\vec{\phi}(y_i, s, a, d) \quad \text{- feature function}$$

16

*Select word labels based on sentence + dependency info*

E.g., *"**Build** **cities** near **rivers** **or** **ocean.**"*

*Word index $j$, sentence $y$, dependency info $q$*

$$p(e_j | j, y, q) \propto e^{\vec{v} \cdot \vec{\psi}(e_j, j, y, q)}$$

*Predicate label $e_j$ = { **action**, **state**, **background** }*

*Log-linear model:*
$$\vec{v} \quad \text{- weight vector}$$
$$\vec{\psi}(e_j, j, y, q) \quad \text{- feature function}$$

17

# Final Q function approximation

*Predict expected value of candidate action*

*State $s$, candidate action $a$*

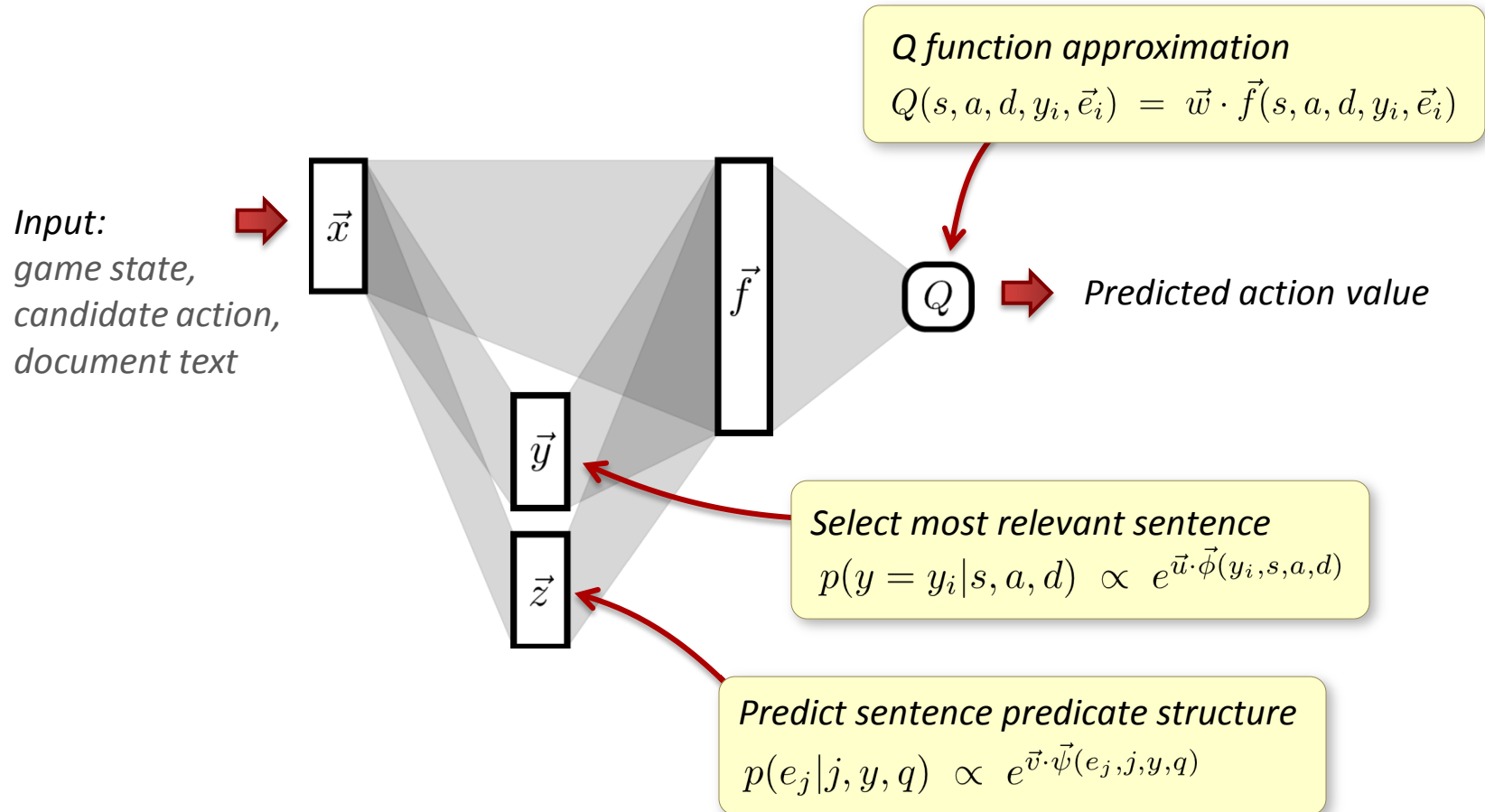$$Q(s, a, d, y_i, \vec{e}_i) = \vec{w} \cdot \vec{f}(s, a, d, y_i, \vec{e}_i)$$

*Document $d$, relevant sentence $y_i$, predicate labeling $\vec{e}_i$*

*Linear model:*
$\begin{cases} \vec{w} & \text{- weight vector} \\ \vec{f}(s, a, d, y_i, \vec{e}_i) & \text{- feature function} \end{cases}$

# Model Representation

Multi-layer neural network: *Each layer represents a different stage of analysis*



**Input:**
*game state,*
*candidate action,*
*document text*

*Q function approximation*
$$Q(s, a, d, y_i, \vec{e}_i) = \vec{w} \cdot \vec{f}(s, a, d, y_i, \vec{e}_i)$$

*Predicted action value*

*Select most relevant sentence*
$$p(y = y_i | s, a, d) \propto e^{\vec{u} \cdot \vec{\phi}(y_i, s, a, d)}$$

*Predict sentence predicate structure*
$$p(e_j | j, y, q) \propto e^{\vec{v} \cdot \vec{\psi}(e_j, j, y, q)}$$

# Parameter Estimation

*Objective*: Minimize *mean square error* between
predicted utility $Q(s, a, d)$
and observed utility $R(s_\tau)$

*Game rollout*



State $s$    Action $a$

Predicted utility: $Q(s, a, d)$          Observed utility: $R(s_\tau)$

# Parameter Estimation

*Method*:  Gradient descent  $-$  i.e.,  Backpropagation.

*Parameter updates*:

$$\vec{u}_i \leftarrow \vec{u}_i + \alpha_u \left[ Q - R(s_\tau) \right] Q \, \vec{x} \left[ 1 - p(y_i | \cdot) \right]$$

$$\vec{v}_i \leftarrow \vec{v}_i + \alpha_v \left[ Q - R(s_\tau) \right] Q \, \vec{x} \left[ 1 - p(e_i | \cdot) \right]$$

$$\vec{w} \; \leftarrow \; \vec{w} + \alpha_w \left[ Q - R(s_\tau) \right] \vec{f}(s, a, d, y_i, z_j)$$

# Features

**State features:**

- *Amount of gold in treasury*

- *Government type*

- *Terrain surrounding current unit*

**Action features:**

- *Unit type (settler, worker, archer, etc)*

- *Unit action type*

**Text features:**

- *Word*

- *Parent word in dependency tree*

- *Word matches text label of unit*

# Experimental Domain

Game:

- *Complex, stochastic turn-based strategy game Civilization II.*
- *Branching factor: $10^{20}$*

Document:

- *Official game manual of Civilization II*

Text Statistics:

*Sentences:* **2083**

*Avg. sentence words:* **16.7**

*Vocabulary:* **3638**

SID MEIER'S
# CIVILIZATION II
THE ULTIMATE VERSION OF THE BEST-SELLING STRATEGY GAME

Preface
&
Instruction Manual

At the start of the game, your civilization consists of a single band of wandering nomads. This is a settlers unit. Although settlers are capable of performing a variety of useful tasks, your first task is to move the settlers unit to a site that is suitable for the construction of your first city. Finding suitable locations in which to build cities, especially your first city, is one of the most important decisions you make in the

# Experimental Setup

Game opponent:

- *Built-in AI of Game.*
- *Domain knowledge rich AI, built to challenge humans.*

Primary evaluation:

- *Games won within first 100 game steps.*
- *Averaged over 200 independent experiments.*
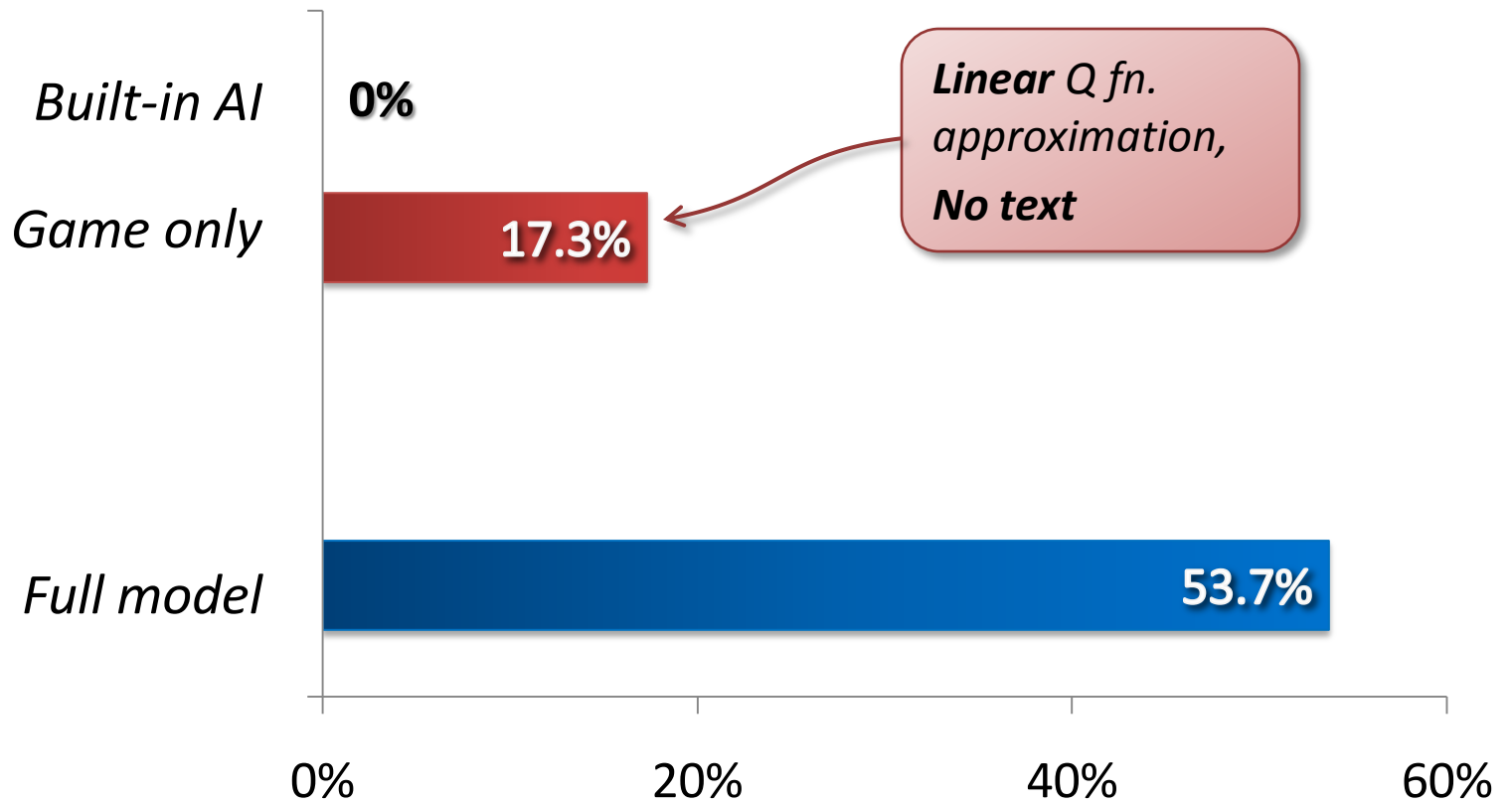- *Avg. experiment runtime: 1.5 hours*

Secondary evaluation:

- *Full games won.*
- *Averaged over 50 independent experiments.*
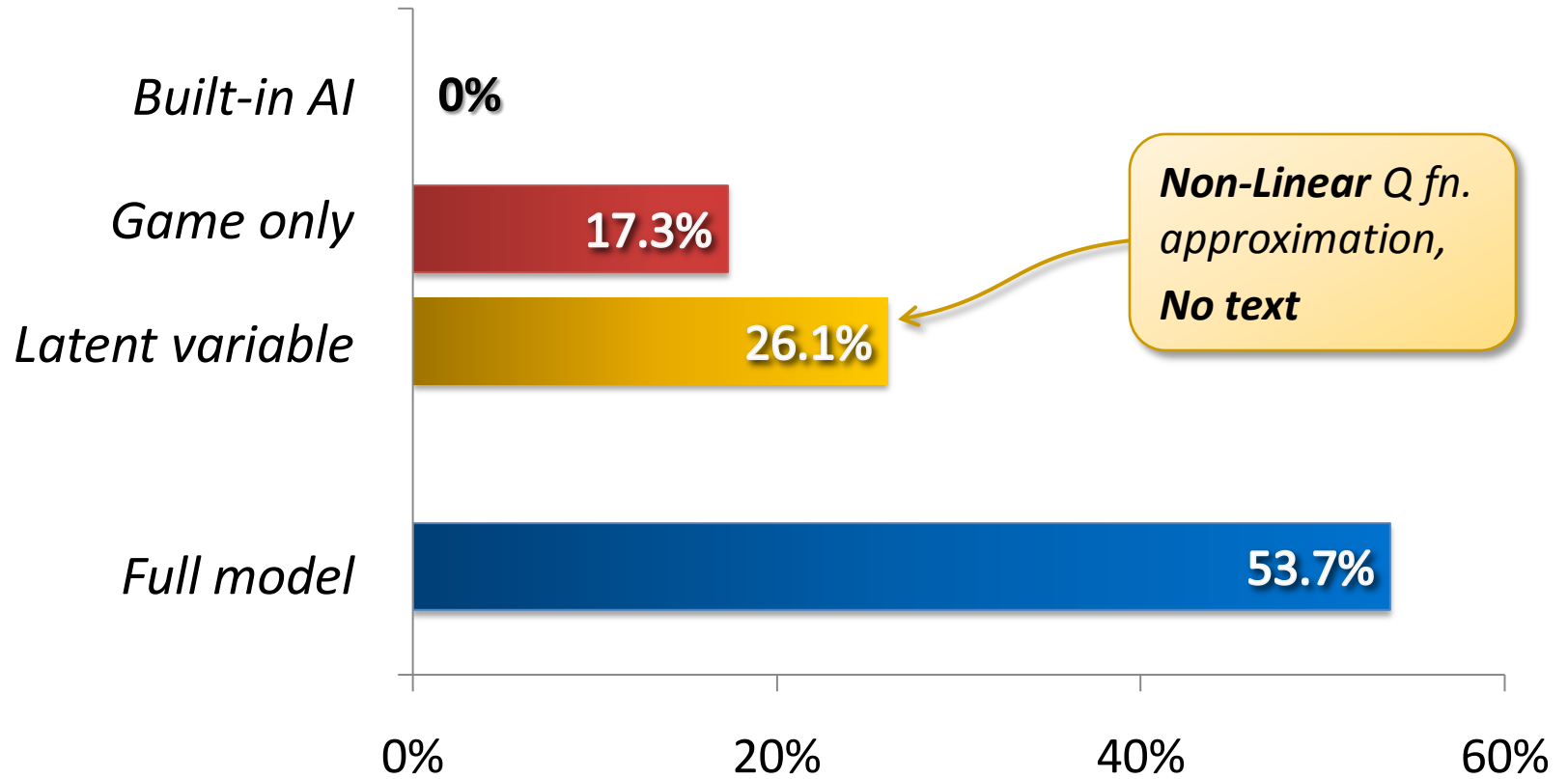- *Avg. experiment runtime: 4 hours*

# Results



Built-in AI    **0%**

Full model    **53.7%**

0%      20%      40%      60%

*% games won in 100 turns, averaged over 200 runs.*

# Does Text Help ?

**Built-in AI**   **0%**

**Game only**   17.3%

> *Linear* Q fn. approximation, **No text**

**Full model**   53.7%

0%    20%    40%    60%

*% games won in 100 turns, averaged over 200 runs.*

# Text vs. Representational Capacity



Built-in AI **0%**

Game only **17.3%**

Latent variable **26.1%**

**Non-Linear** Q fn. approximation, **No text**

Full model **53.7%**

0%    20%    40%    60%

*% games won in 100 turns, averaged over 200 runs.*

# Linguistic Complexity vs. Performance Gain



*% games won in 100 turns, averaged over 200 runs.*

# Results: Sentence Relevance

**Problem:** *Sentence relevance depends on game state.*

*States are game specific, and not known a priori!*

**Solution:** *Add known non-relevant sentences to text.*

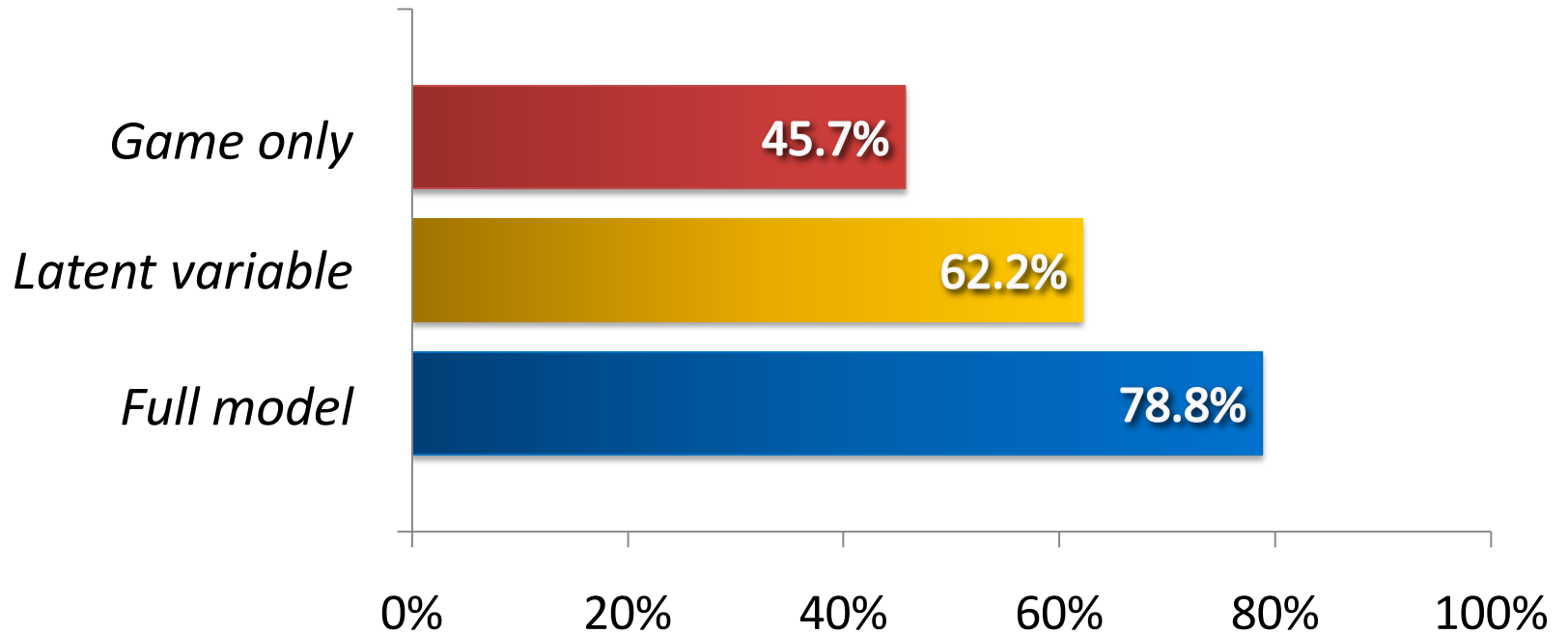*E.g., sentences from the Wall Street Journal corpus.*

**Results:** ***71.8%*** *sentence relevance accuracy…*

*Surprisingly poor accuracy given game win rate!*

# Results: Sentence Relevance

# Results: Full Games



*Percentage games won, averaged over 50 runs*

# Related Work

*Grounded Language Acquisition: Instruction Interpretation*

*Branavan et al. 2009, 2010,   Vogel & Jurafsky 2010*

- Imperative descriptions of action sequences
- Assume relevance of text to current world state

*Language Analysis in Games*

*Eisenstein et al. 2009*

- Extract high-level semantic representation from text
- Learn game rules from labeled traces  +  extracted formulae

*Gorniak & Roy 2005*

- Interpret spoken commands to control game character
- Learn from labeled parallel corpus

# Conclusions

- Human knowledge encoded in natural language can be automatically leveraged to improve control applications.

- Environment feedback is a powerful supervision signal for language analysis.

- Method is applicable to control applications that have an inherent success signal, and can be simulated.

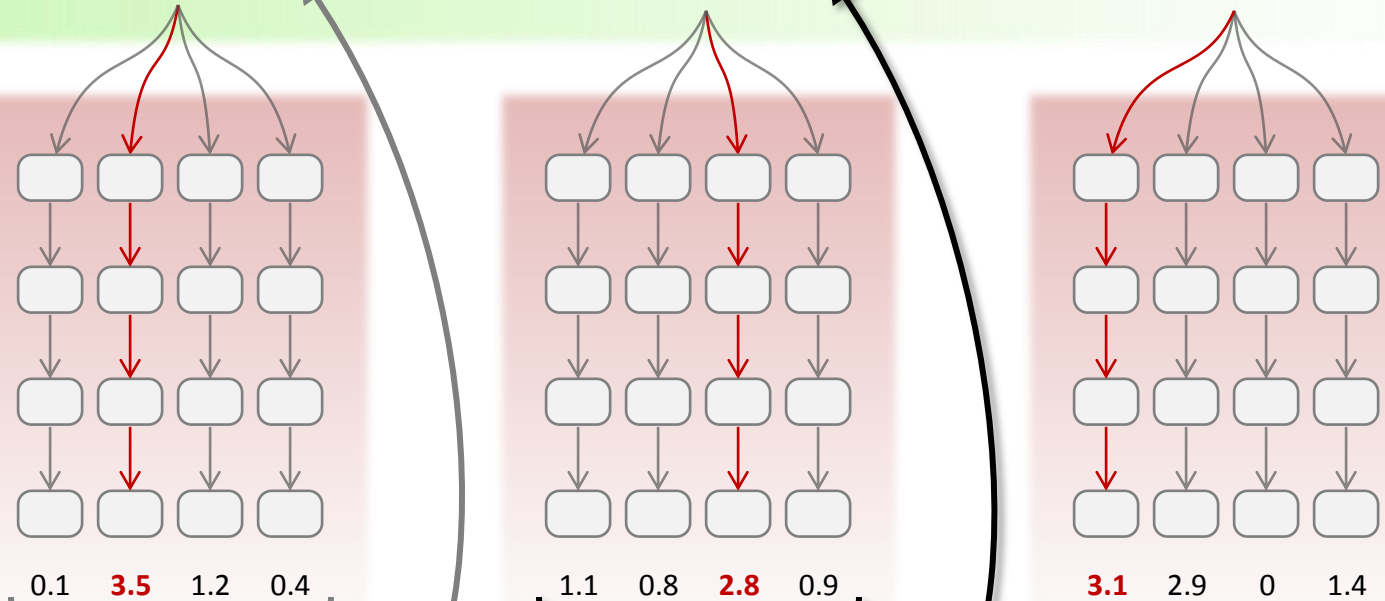*Code, data & experimental framework available at:*
http://groups.csail.mit.edu/rbg/code/civ

# Monte-Carlo Search: Summary



Game states and actions

state 1    action 1    state 2    action 2    state 3

*Monte-Carlo Rollouts (simulations)*

0.1  **3.5**  1.2  0.4        1.1  0.8  **2.8**  0.9        **3.1**  2.9  0  1.4
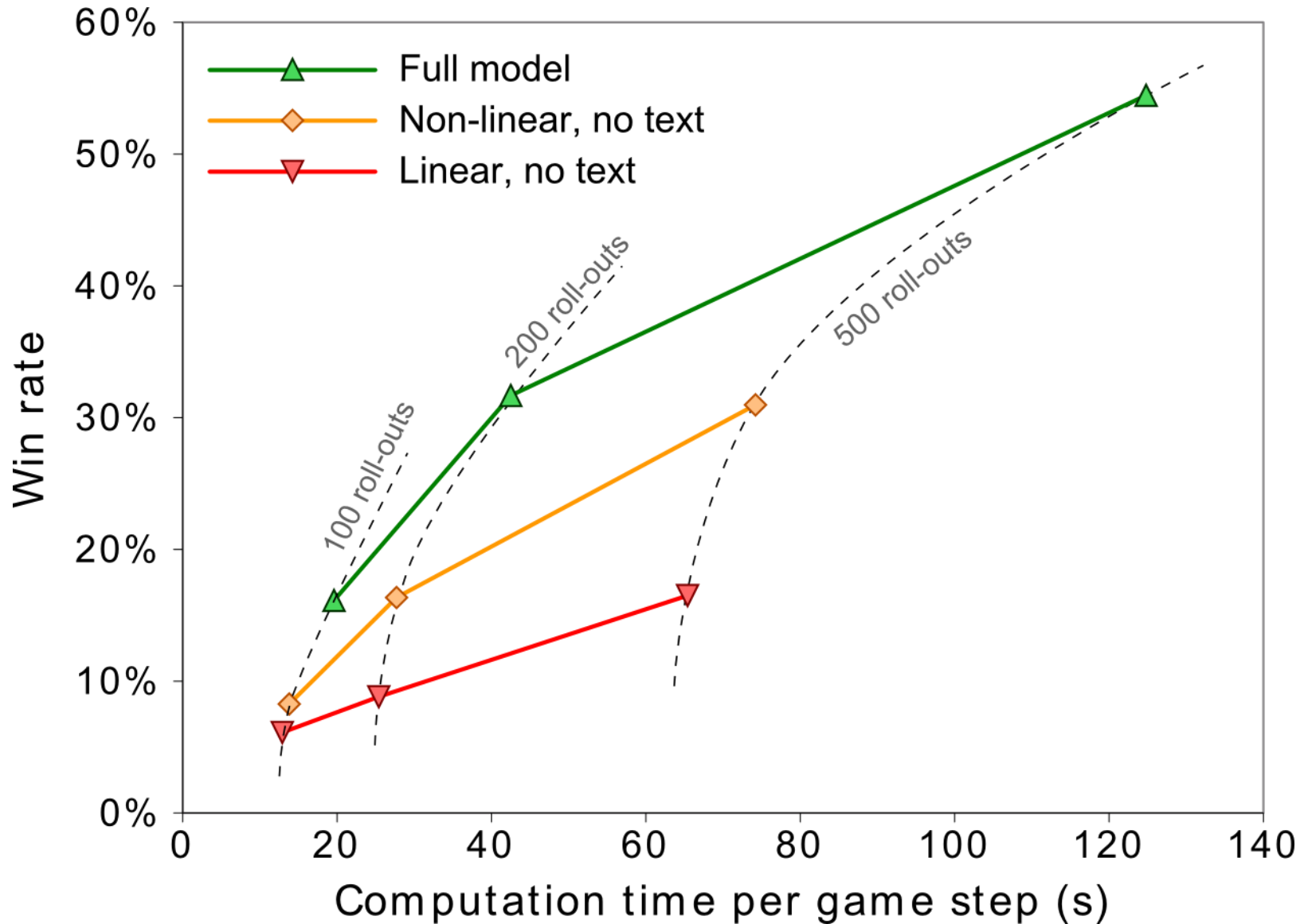
*Use observed rollout scores to select game action*

# Model Complexity, Time and Performance

# Dependency Information