# Human Priors in Hierarchical Program Induction

Mark K Ho*, Sophia Sanborn*, Frederick Callaway*, David Bourgin, and Thomas L. Griffiths
**University of California, Berkeley**

mark_ho@berkeley.edu, sanborn@berkely.edu, fredcallaway@berkeley.edu

## Background

► Human problem-solving behavior is organized into **rich hierarchical structure** [1, 2].
► Inferring this structure is essential for interpreting behavior and anticipating how others will act.

## Research Goals

► Previous work cast hierarchy learning as an efficient coding problem and found people often generate shorter programs to solve problems [3].
► Here, we examine alternative **program features** that constrain how people **interpret** hierarachically organized behavior.

## Programs as Problem Solutions

► A program ($\pi$) is a set of subprocesses, $\sigma_i$, which are sequences of **primitive actions** or **calls to other subprocesses**.
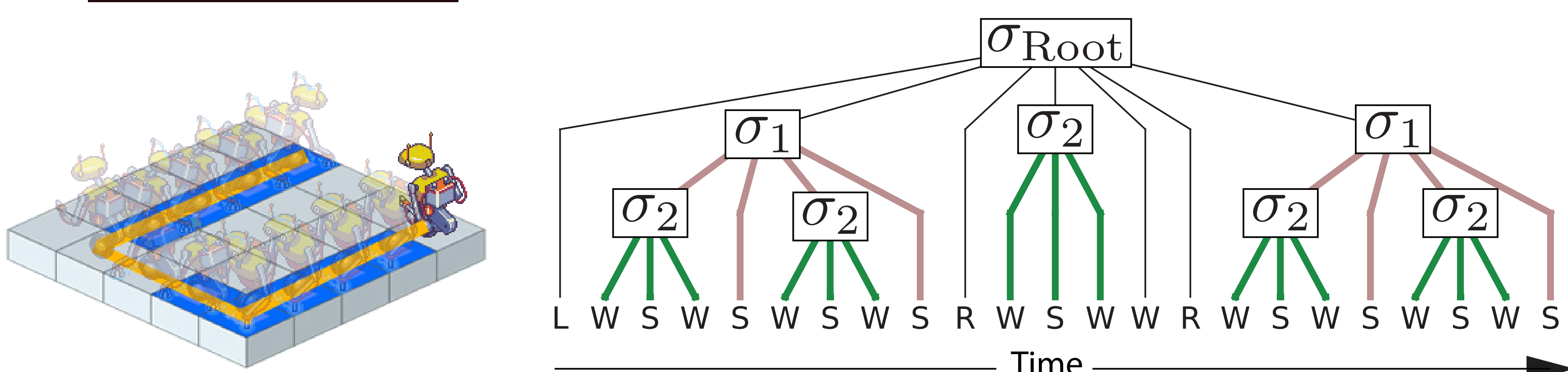
$$\pi = \{\sigma_{\text{Root}}, \sigma_1, \sigma_2\}$$
$$\sigma_{\text{Root}} = (\texttt{Left}, \sigma_1, \texttt{Right}, \sigma_2, \texttt{Walk}, \texttt{Right}, \sigma_1)$$
$$\sigma_1 = (\sigma_2, \texttt{Switch}, \sigma_2, \texttt{Switch})$$
$$\sigma_2 = (\texttt{Walk}, \texttt{Switch}, \texttt{Walk})$$

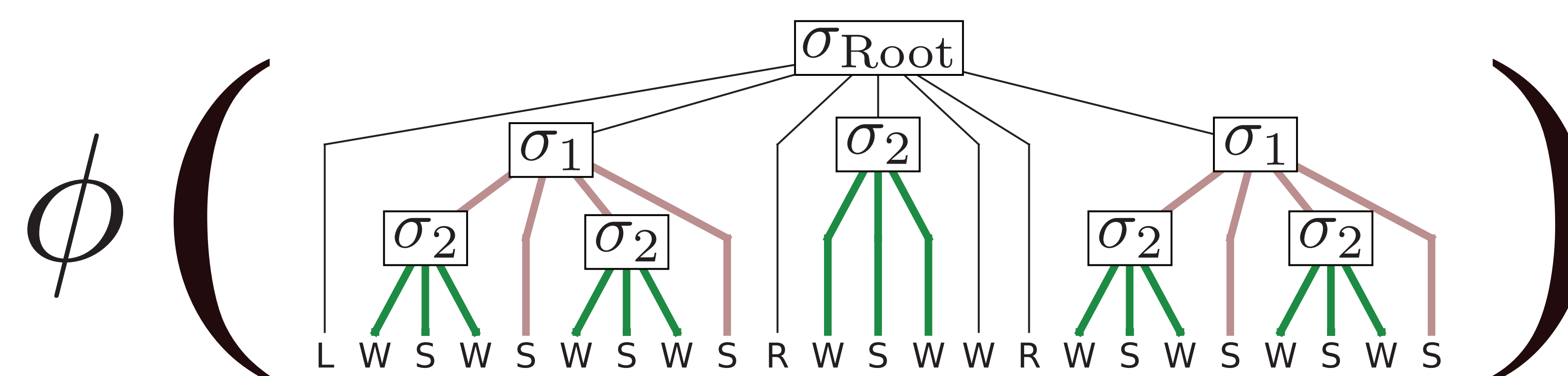► Executing a program produces a **state-action trace** as well as an **execution tree**.



## Inferring Programs

► Inducing a program from a trace can be expressed as probabilistic inference:

$$p(\pi \mid \zeta) = \frac{p(\zeta \mid \pi)p(\pi)}{\sum_{\pi'} p(\zeta \mid \pi')p(\pi')}$$

► The program prior is a function of weighted program features.

$$p(\pi; \theta) \propto \exp\{\theta^\top \phi(\pi)\}$$



$$=
\begin{bmatrix}
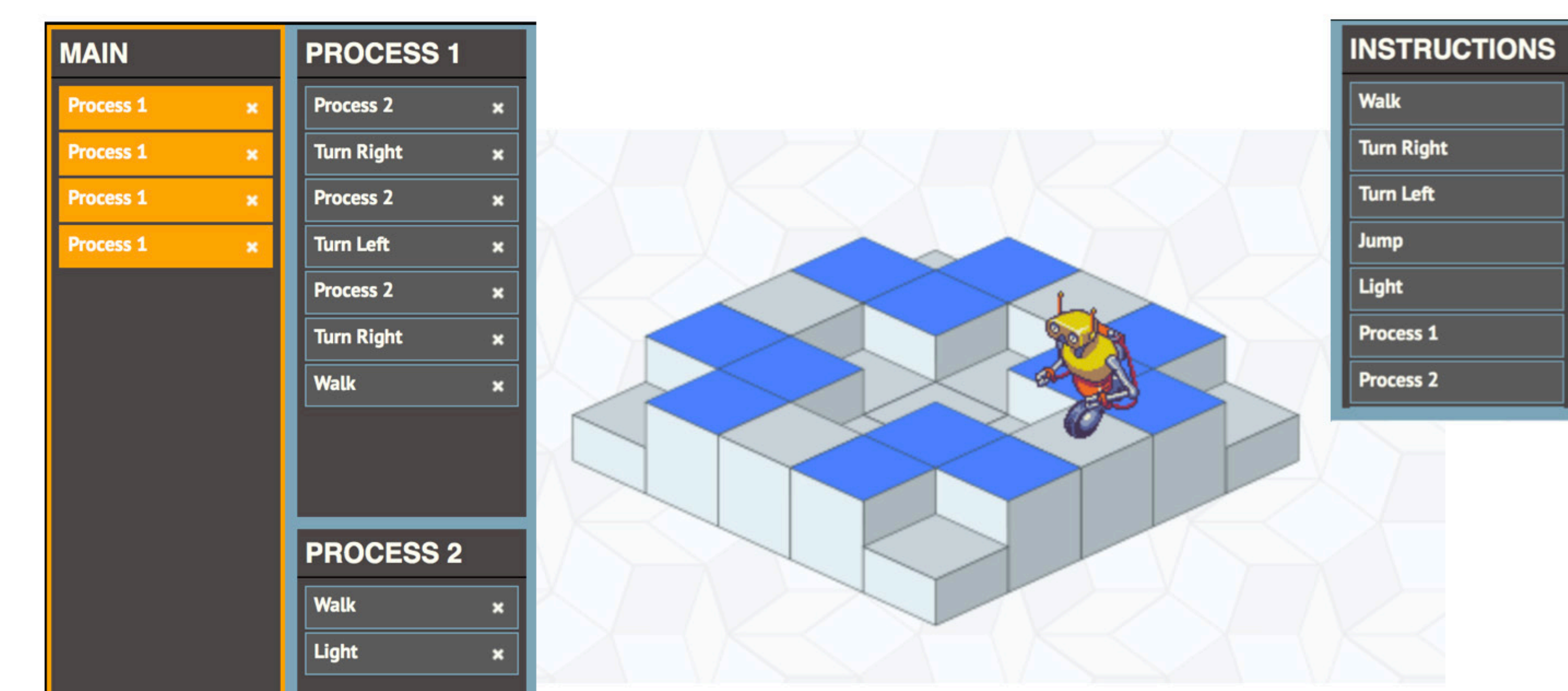14 \\ 3 \\ 2 \\ 1.7 \\ 0.66 \\ 0.81 \\ 0.46 \\ 7 \\ 1.67 \\ 1.5 \\ 1.38 \\ 0.58
\end{bmatrix}
\begin{array}{l}
\leftarrow \text{Total Program Length} \\
\leftarrow \text{Execution Tree Depth} \\
\leftarrow \text{Number of Subprocesses} \\
\leftarrow \text{Subprocess Length (S.D.)} \\
\leftarrow \text{Action/Subprocess Entropy (M)} \\
\leftarrow \text{Action-Call Entropy (M)} \\
\leftarrow \text{Process-Call Entropy (M)} \\
\leftarrow \text{Root Length} \\
\leftarrow \text{Children per Subprocess (M)} \\
\leftarrow \text{Parents per Subprocess (M)} \\
\leftarrow \text{Subprocess Entropy (M)} \\
\leftarrow \text{Subprocess to Action Ratio (M)}
\end{array}$$

► Given features $\phi$, trace $\zeta$, and program $\pi$, we want to estimate the feature weights $\theta$
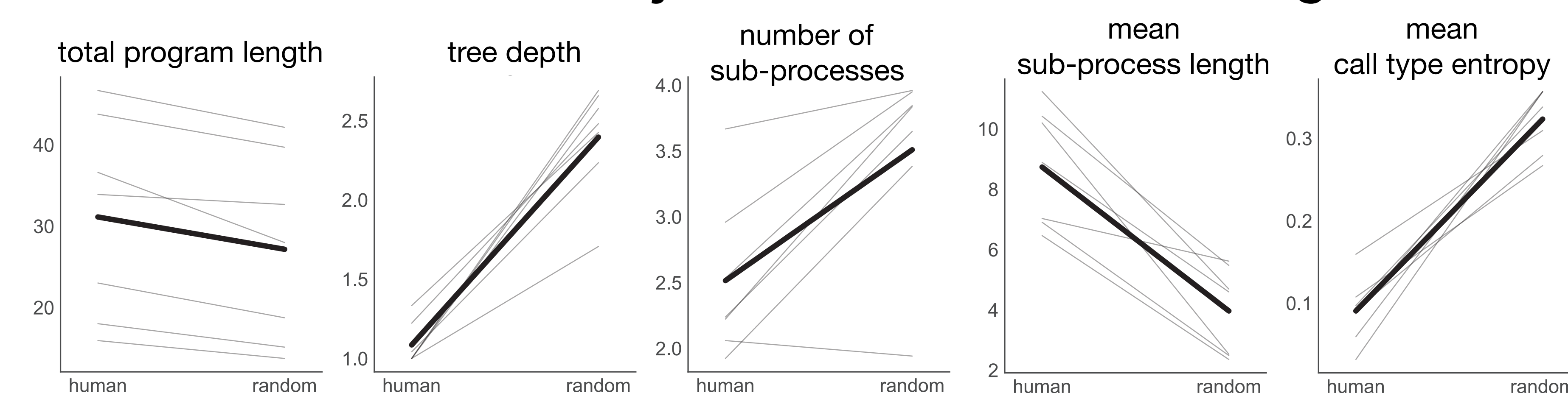
$$p(\theta \mid \pi, \zeta) \propto \frac{p(\zeta \mid \pi)p(\pi; \theta)}{p(\zeta; \theta)}p(\theta)$$

## Estimating Human Priors

► 77 participants observed human-generated solution traces in the Lightbot domain [1] and tried to reconstruct the original programs that generated the traces.



### Empirical Feature Weights for Humans vs. Randomly Generated Consistent Programs



## Conclusions

► Participants prefer programs that have shallower trees and use fewer, longer subprocesses.
► These findings may reflect working memory constraints that limit the structural complexity of induced programs.
► In ongoing work, we are investigating the relationships between these features and their individual contributions to program induction, program generation, and perceived program complexity.

**References**
[1] Simon, H. A. (1991) The Architecture of Complexity.
[2] Solway, A., Diuk, C., Cordova, N., Yee, D., Barto, A., Niv, Y., & Botvinick, M. (in press). Optimalbehavioral hierarchy. PLOS Computational Biology.
[3] Sanborn, S., Bourgin, D., Chang, M., & Griffiths, T. (2018). Representational efficiency outweighs action efficiency in human program induction. In Proceedings of the 40th annual cognitive science society.