

DEEP GAUSSIAN PROCESSES FOR LARGE DATASETS

JAMES HENSMAN^{a,b}, ANDREAS DAMIANOU^{a,b} AND NEIL D. LAWRENCE^{a,b}

PROBLEM

We wish to develop a Deep Belief Network where the transformation between layers is probabilistic and modelled with Gaussian processes.

- The function $f_i(\cdot)$ is modelled with a GP
- Assume noise β_i at each level
- The mapping instantiations $\bar{f}_i = f_i(h_i)$ can be marginalised out analytically: $p(h_i|h_{i-1}) =$

../diagrams/graphical_collapsed.pdf

$$= \int p(h_i|\bar{f}_i)p(\bar{f}_i|h_{i-1}) d\bar{f}_i$$

$$= \mathcal{N}(h_i|0, K(h_{i-1}, h_{i-1}) + \beta_i I)$$

- How to learn the intermediate hidden layers?
- How to efficiently train the model?

INDUCING VARIABLES

- Also marginalise out hidden spaces: learn a posterior $q(h_i)$ for each layer.
- Use inducing points z_i : $u_i = f_i(z_i)$.
- Let $q(u_i) = \mathcal{N}(u_i|m_i, \Sigma_i)$.
- Inducing points play the same role as $\{h_i, \bar{f}_i\}$ pairs but result in low rank representations of the covariance matrices.

../diagrams/graphical_collapsed.pdf

- Inducing points are variational parameters allowing us to lower bound the evidence: $\mathcal{F} \leq \log p(y)$

Given the above, how can we define $p(h_i|h_{i-1}, u_i)$?

INFERENCE STRATEGIES

We need to deal with $q(h)$ and $q(u)$. Three strategies:

1. Collapse out u (Damianou et al., AISTATS 2013)
 - Optimize $q(h)$
2. Maintain $q(h)$ and $q(u)$
 - EM-style optimisation for $q(u)$ and $q(h)$
3. Compress $q(h)$ into $q(u)$ using $p(h|u)$

VARIATIONAL COMPRESSION

Consider one layer of our inference problem. Use the conditional distribution as the variational distribution:

$$p(y|u) = \frac{p(y|h)p(h|u)}{p(h|y)}$$

$$\log p(y|u) = \mathbb{E}_{p(h|u)} [\log p(y|h)] + \text{KL}[p(h|u) \parallel p(h|y)]$$

Small if u explains h very well

We can compute the marginal distribution in the variational approximation easily:

$$q(h_i) = \int \underbrace{p(h_i|u_i)}_{\sim \mathcal{N}} \underbrace{q(u_i)}_{\sim \mathcal{N}} du_i$$

../diagrams/graphical_collapsed.pdf

- Given X and a fixed $q(u_1)$, we can compute $q(h_1)$
- For a fixed $q(h_1)$, we can variationally propagate using $q(u_2)$ to get $q(h_2)$ (blue arrows)
- Continue to feed-forward to the bottom layer. The variational propagation at each layer introduces a penalty (regularizing) term which affects the bound on the marginal likelihood
- Applying the chain-rule leads to backpropagation (red arrows), but with *Gaussian* messages passed layer-to-layer

INFERENCE FOR LARGE DATASETS

How can we handle large datasets?

Stochastic variational inference (SVI):

SVIGP-style: [Hensman et al., UAI 2013]

- Represent the parameters $\{m, \Sigma\}$ of $q(u)$ in two equivalent ways:
 - Canonical form: $\theta = \{\Sigma^{-1}m, -\frac{1}{2}\Sigma^{-1}\}$
 - Expectation form: $\eta = \{m, mm^T + \Sigma\}$
- Treat u as *global variables*. This allows for the factorisation of the contributions of every input/output pair $\{x_{l,i}, y_i\}$.
- Optimise the parameters using the natural gradients of $q(u)$:

$$\theta_{t+1} = \theta_t + s \frac{\mathcal{F}}{d\eta},$$

where s is the learning step

Adapting the learning step s :

- Stochastic optimisation is very sensitive to s .
- [Ranganath et al. ICML 2013] dynamically adapt s to minimize the expected loss between the parameter vector after the stochastic variational update, θ_{t+1} , and the vector after a full variational update θ_{t+1}^*
- Here we consider the loss in the KL sense, considering the involved distributions:

$$\text{KL}[q(u|\theta^*) \parallel q(u|\theta)]$$

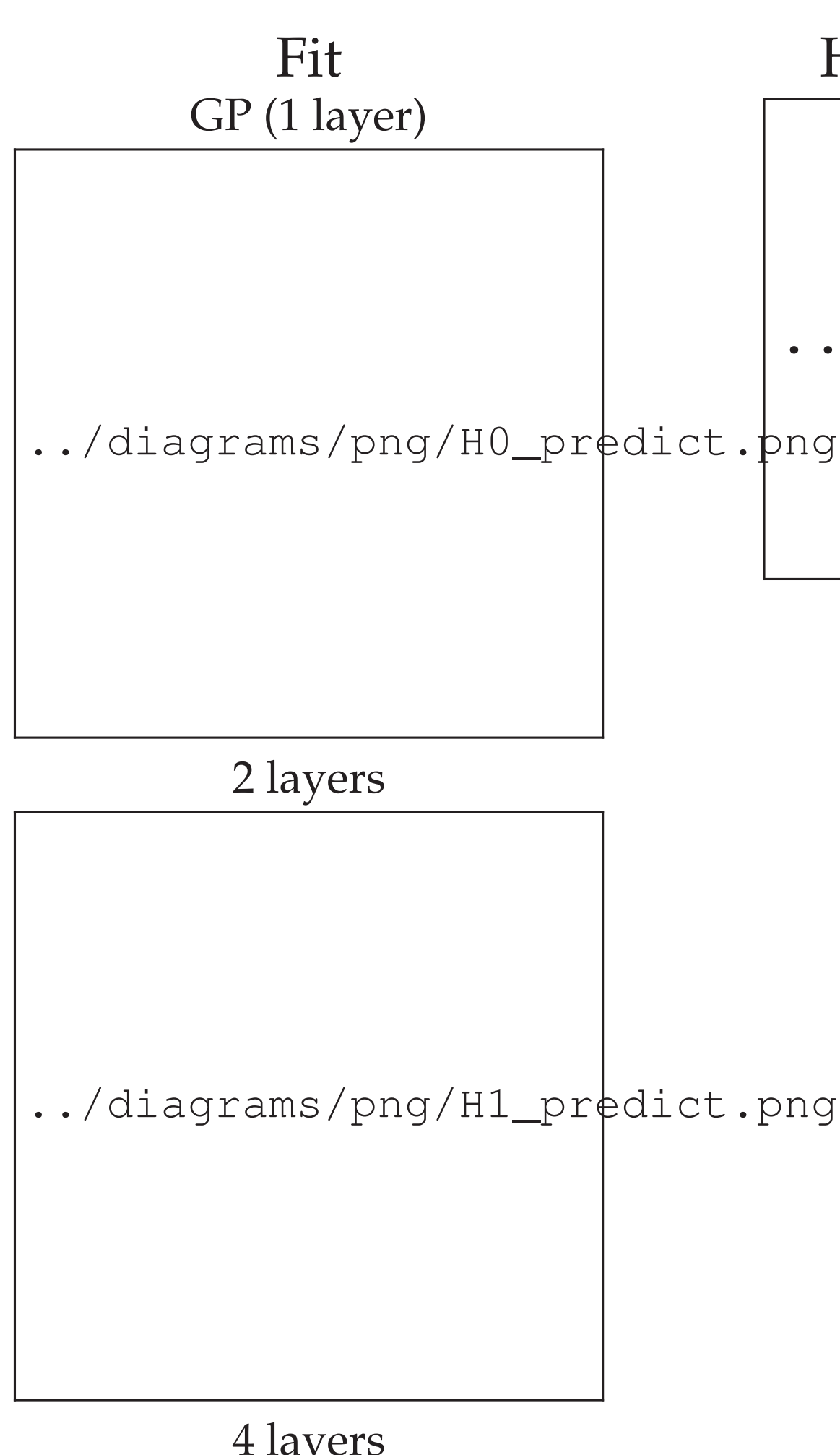
This takes into account the geometry of the parameter space.

TODO

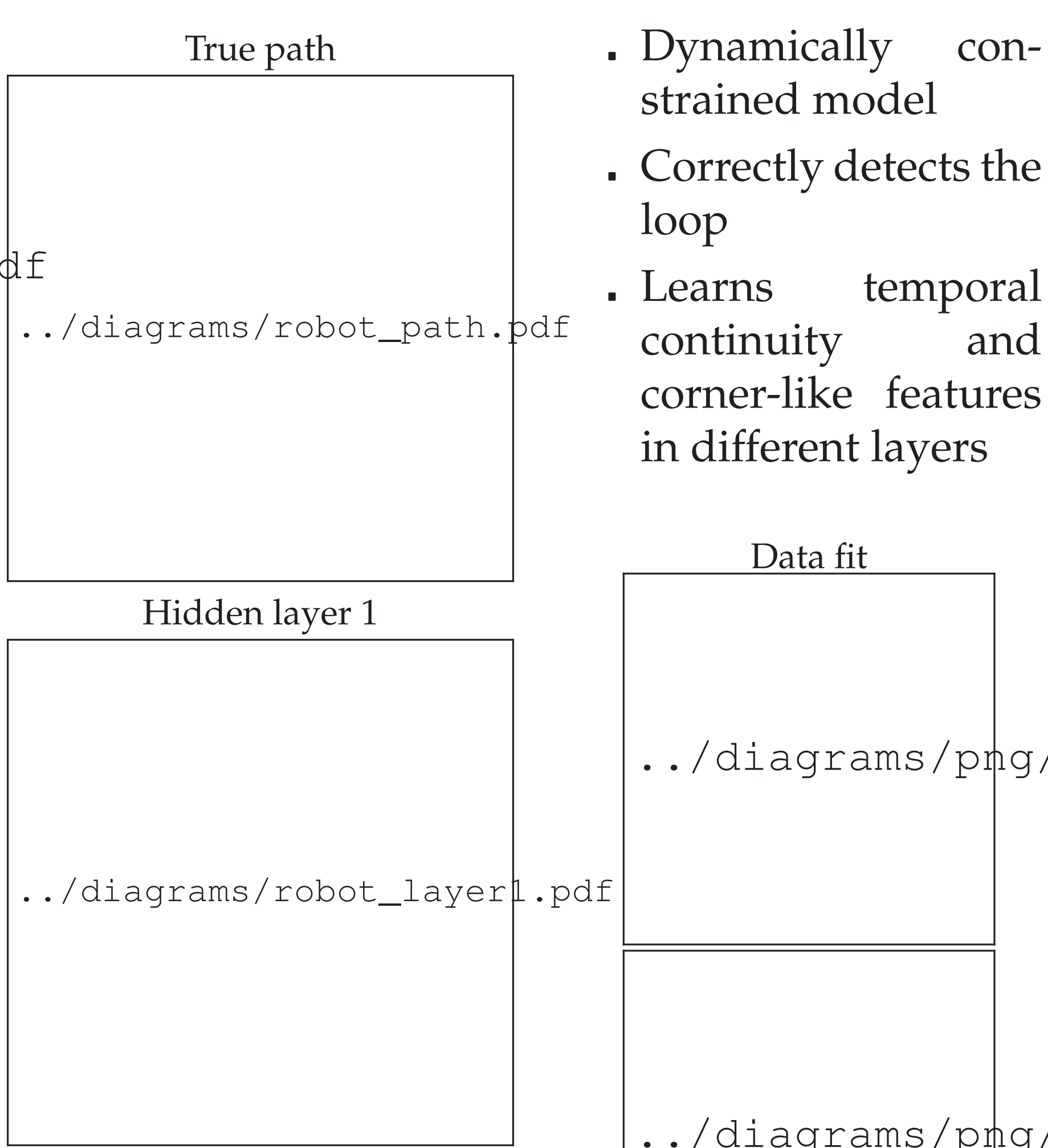
- Currently, SVI inference is implemented only for 1-layer models
- Extend SVI inference framework in deep models
- Training scheme combining optimisation of variational and kernel parameters
- Fix initialisation issues
- Explore auto-encoder architectures

EXPERIMENTS

Toy problem

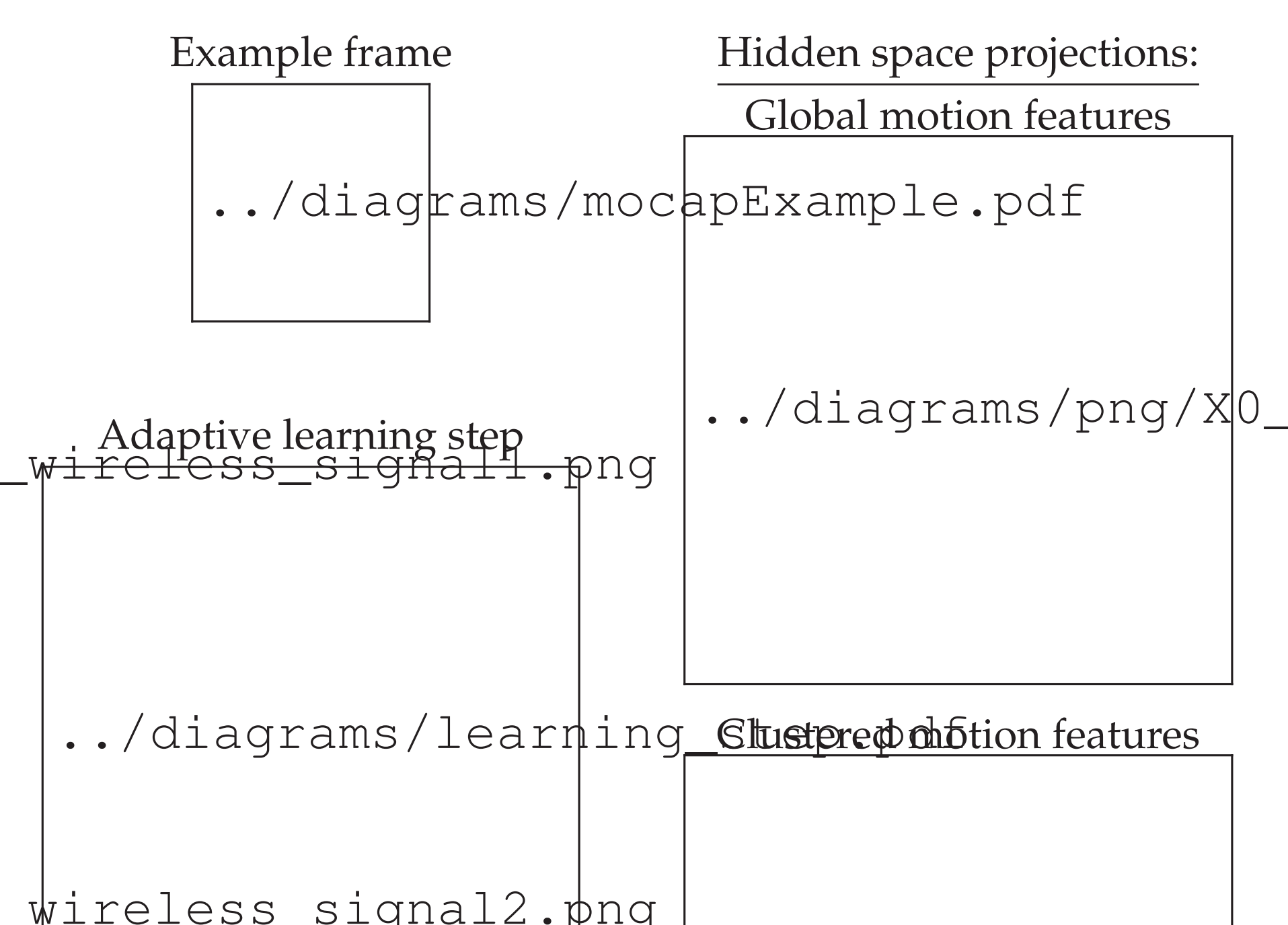


Loop detection in robotics



Big Data

- 12 Subjects, 95 diverse motions, 20K datapoints
- Learns a general model of human motion
- Outperforms Bayesian GP-LVM (trained on subsets) for reconstructing part of test body parts
- We considered a 1-layer model but used SVI inference with adaptive learning step



^a University of Sheffield, Department of Computer Science, UK
 Corresponding author: james.hensman@sheffield.ac.uk

^b Sheffield Institute for Translational Neuroscience (SITRN), UK
 Example data fits for 2 of the 30 output dimensions
 ../diagrams/robot_layer2.pdf