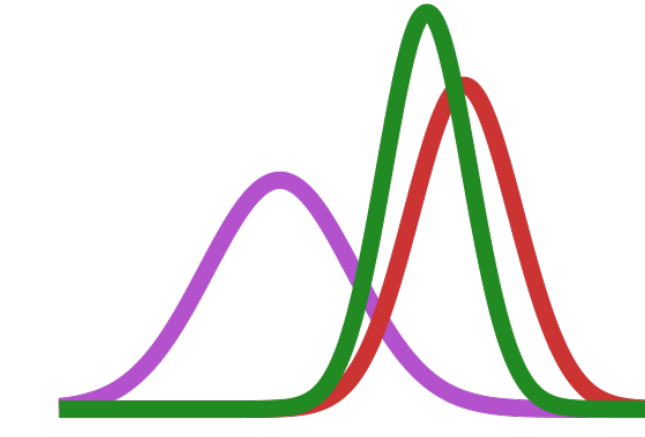


DynamicPPL: Stan-like Speed for Dynamic Probabilistic Models

Mohamed Tarek^{1,2}, Kai Xu³, Martin Trapp⁴, Hong Ge⁵, and Zoubin Ghahramani^{5,6}

¹University of New South Wales at Canberra; ²Pumas-AI Inc.; ³ University of Edinburgh; ⁴ Graz University of Technology; ⁵ University of Cambridge; ⁶ Google Research

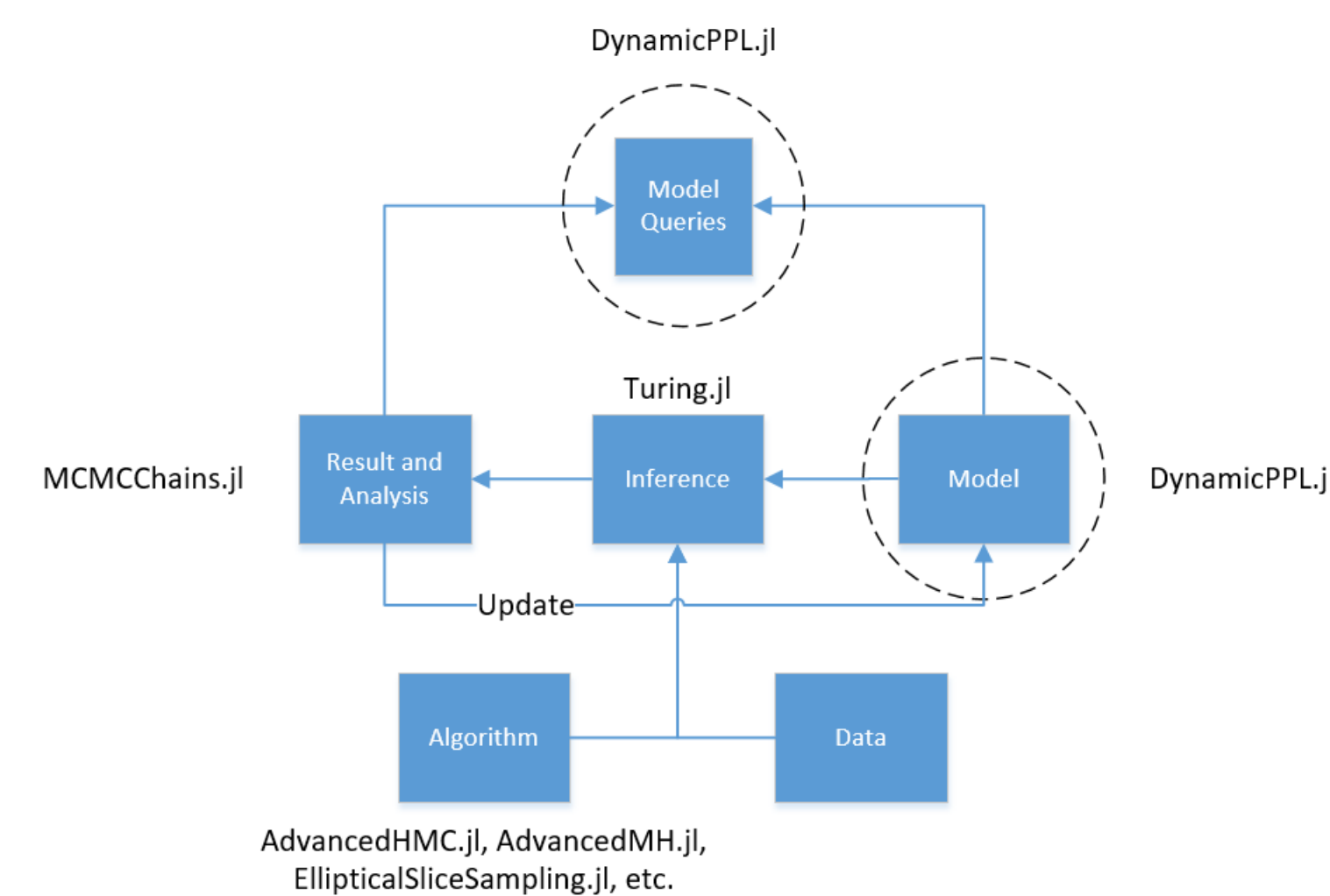


1. Introduction

Practical Bayesian inference for probabilistic models with dynamic dimensionality and variable types is a challenging problem in probabilistic programming languages (PPLs). General-purpose PPLs enable users to introduce new model parameters anywhere in a model definition, freely influence control flows with random variables, and often eliminate the need to specify variable types. Such flexibility often leads to a more productive probabilistic modelling workflow – models are easy to read and write. Such increased productivity is critical for a broader adoption of Bayesian modelling and approximate inference. In this poster, we present DynamicPPL which enabled a significant speedup of inference in Turing – a general-purpose PPL implemented in the Julia programming language. Our approach leverages incremental tracing and type specialisation. We show through an extensive set of benchmarks that the proposed approach significantly accelerates inference in dynamic models (often an order of magnitude faster than a previous version of Turing) and achieves state-of-the-art performance for static models that is close to and sometimes better than Stan’s (which is implemented in C++).

The main contribution of this poster is the **DynamicPPL** Julia package, a user-friendly and modular PPL frontend that makes use of incremental tracing and type specialization of model parameters in probabilistic programs allowing **Turing** [2] to achieve speeds close to and sometimes faster than Stan even for static models while efficiently supporting a family of dynamic models.

2. Overview of TuringLang Organisation

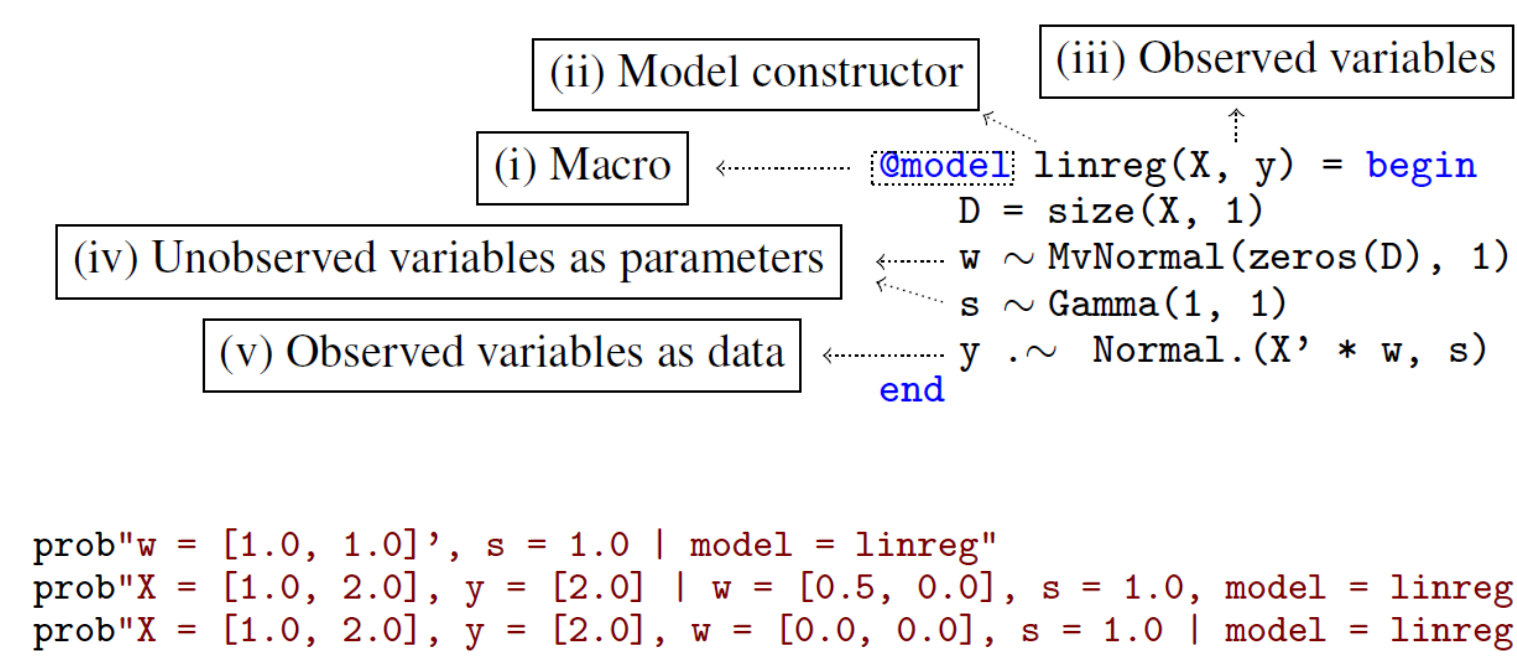


3. Dynamic vs Static Typing

- Dynamic PPLs lack the random variables’ type information at compile time.
- Type information of variables is paramount to many compiler optimisations required to generate efficient machine code; Python is slower than C.
- The Julia programming language [1] is a fast dynamically typed programming language; its design and compiler allow it to behave both as a dynamic programming language as well as a static one in different contexts, thus combining the advantages of both paradigms.

4. Design Principles

- User-friendly syntax

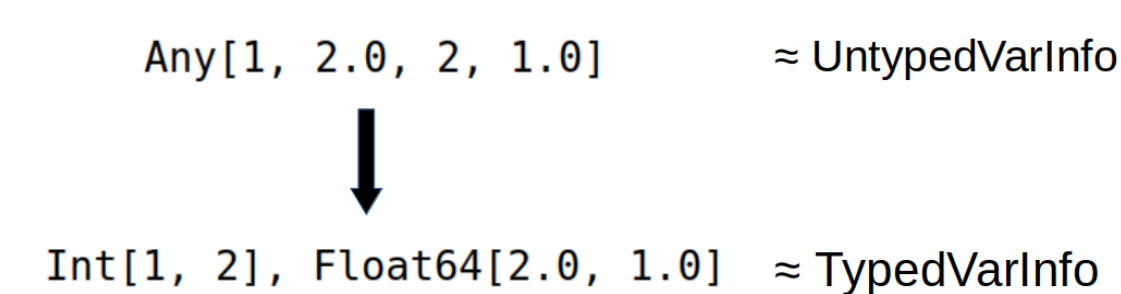


- Inter-operability with the Julia ecosystem
 - Distributions.jl for [distributions](#)
 - ForwardDiff.jl, Tracker.jl, ReverseDiff.jl and Zygote.jl for [automatic differentiation](#)
 - Flux.jl for [neural networks](#)
 - Optim.jl for [optimisation](#) used in maximum a-posteriori and maximum likelihood estimation
 - CUDA.jl for [GPU data parallelism](#)
 - DifferentialEquations.jl for [differential equation solvers](#)
 - Memoization.jl for [memoising expensive functions](#) in dynamic Gibbs sampling
 - Arbitrary Julia code allowed including [custom distributions](#)
- Performance: comparable performance to Stan

5. Trace Type Specialisation

Trace type specialisation in **DynamicPPL + Turing**:

1. Run the model using **UntypedVarInfo**
2. Specialise the container types in the trace creating a **TypedVarInfo**
3. Run the rest of inference using a trace specialised for the model
4. (WIP) Incrementally “expand” the trace type as needed using a **MixedVarInfo**



The use of **TypedVarInfo** allows for optimised code generation for:

- Static models, and
- A class of dynamic models where the random variable containers’ names and types are static, but the number of random variables and their distributions can be dynamic.

The **MixedVarInfo** work will enable efficient code generation of fully dynamic models, where type information can also be dynamic, with 0 overhead for the currently efficient classes of models. New type information can get exploited incrementally to re-generate efficient code using the Julia compiler.

6. Benchmarks

- Static HMC benchmarks of Turing against Stan using AdvancedHMC.jl + ReverseDiff.jl from Turing.jl. Time is in seconds.

	Forward evaluation		Gradient evaluation	
	Stan	Turing	Stan	Turing
Gaussian unknown	3.296×10^{-5}	1.898×10^{-5}	4.23×10^{-5}	1.25×10^{-4}
Hierarchical Poisson	1.049×10^{-5}	7.067×10^{-6}	2.44×10^{-5}	7.84×10^{-5}
High dimensional Gaussian	1.307×10^{-3}	4.579×10^{-5}	1.81×10^{-3}	1.44×10^{-4}
Semi-supervised HMM	6.374×10^{-4}	2.939×10^{-4}	5.97×10^{-4}	3.81×10^{-3}
Latent Dirichlet allocation	1.624×10^{-3}	7.583×10^{-5}	1.89×10^{-3}	1.62×10^{-3}
Logistic regression	1.137×10^{-3}	1.525×10^{-3}	1.59×10^{-3}	3.23×10^{-3}
Naïve Bayes	3.672×10^{-4}	1.246×10^{-5}	3.23×10^{-4}	3.60×10^{-5}
Stochastic volatility	1.298×10^{-4}	7.553×10^{-4}	1.40×10^{-4}	7.90×10^{-4}

- Static HMC benchmarks of Turing against Stan using AdvancedHMC.jl + ReverseDiff.jl from Turing.jl. Time is in seconds.

	Stan	Turing
Gaussian unknown	1.415 ± 0.009	10.417 ± 0.392
Hierarchical Poisson	0.310 ± 0.007	5.223 ± 0.645
High dimensional Gaussian	40.418 ± 2.012	19.576 ± 6.448
Semi-supervised HMM	42.464 ± 0.095	342.516 ± 124.477
Latent Dirichlet allocation	147.210 ± 7.871	91.551 ± 25.194
Logistic regression	105.097 ± 8.748	390.126 ± 8.345
Naïve Bayes	21.282 ± 1.497	4.273 ± 0.158
Stochastic volatility	5.552 ± 0.041	230.465 ± 29.179

7. Conclusion

1. Incremental tracing and type specialization facilitates the acceleration of inference in dynamic PPLs such as **Turing**
2. **DynamicPPL** is a modular, high-performance, dynamic PPL frontend implementation with a user-friendly and extensible interface
3. In the future, we hope to extend **DynamicPPL** and **Turing** to support nested modeling, BUGS-style Gibbs sampling and Infer.net-style message passing.

8. References

- [1] Jeff Bezanson et al. “Julia: A fresh approach to numerical computing”. In: *SIAM review* 59.1 (2017), pp. 65–98.
- [2] Hong Ge, Kai Xu, and Zoubin Ghahramani. “Turing: A Language for Flexible Probabilistic Inference”. In: *Proceedings of the Twenty-First International Conference on Artificial Intelligence and Statistics*. Ed. by Amos Storkey and Fernando Perez-Cruz. Vol. 84. Proceedings of Machine Learning Research. Playa Blanca, Lanzarote, Canary Islands: PMLR, 2018, pp. 1682–1690.