

Contents

Notation	xi
1 Introduction	1
1.1 Related Work	3
2 Background	5
2.1 Bayesian Inference and MCMC methods	5
2.2 Probabilistic Programming	10
2.3 Compilation and Metaprogramming in Julia	14
2.4 Automatic Differentiation and Computation Graphs	21
3 Implementation of Dynamic Graph Tracking in Julia	25
3.1 Extended Wengert Lists	26
3.2 Automatic Graph Tracking	27
3.3 Evaluation	31
4 Graph Tracking in Probabilistic Models	33
4.1 Dependency Analysis in Dynamic Models	33
4.2 JAGS-Style Automatic Calculation of Gibbs Conditionals	33
4.3 Evaluation	33
5 Discussion	35
5.1 Future Work	35
List of Algorithms	37

2 Background

This chapter provides the background for the concepts used later in chapters 3 and 4. Initially, it gives a quick overview of Bayesian inference and probabilistic programming in general, necessary to understand the requirements and usual approaches of probabilistic programming systems.

Consequently, the machinery and language used to develop the graph tracking system forming the main part of the work are described. This consists firstly of the basic notions and techniques of the Julia compilation process as well as the language’s metaprogramming capabilities are described, which form the basis of the implementation. Secondly, a short introduction to graph tracking and source-to-source automatic differentiation is given, which contains many ideas and terminology that will be used later, and often provided inspiration.

2.1 BAYESIAN INFERENCE AND MCMC METHODS

Generative modelling is an approach for modelling phenomena based on the assumption that observables can be fully described through some stochastic process. When we assume this process to belong to a specified family of processes, the estimation of the “best” process is a form of learning: if we have a good description of how observations are generated, we can make summary statements about the whole population (descriptive statistics) or predictions about new observations. When observations come in pairs of independent and dependent variables, learning the conditional model of one given the other solves a regression or classification problem.

Within a Bayesian statistical framework, we assume that the family of processes used is specified by random variables related through conditional distributions with densities, which describe how the observables would be generated: some *unobserved variables* are generated from *prior distributions*, and the *observed data* are generated conditionally on the unobserved variables. The goal is to learn the *posterior distribution* of the parameters given the observations, which is a sort of “inverse” of how the problem is specified.

As an example, consider image classification: if we assume that certain percentages of an image data set picture cats and dogs, respectively, the distribution of these labels forms the prior. Given the information which kind of animal is depicted on it, an image can then be generated as a matrix of pixels based on a distribution of images conditioned on labels. The posterior distribution is then conditional distribution of

the label given an image. When we have this information, we can, for example, build a Bayesian classifier, by returning for a newly observed image that label which has the highest probability under the posterior.

This kind of learning is called Bayesian inference since, in the form of densities, the form of the model can be expressed using Bayes’ theorem as the conditional distribution with density¹

$$\overbrace{p(\theta | x)}^{\text{posterior}} = \frac{\overbrace{p(x | \theta)}^{\text{likelihood}} \overbrace{p(\theta)}^{\text{prior}}}{p(x)}, \quad (2.1)$$

where x are the observed data, and θ are the unobserved parameters. The posterior represents the distribution of the unobserved variables as a combination of the prior belief updated by what has been observed (Congdon 2006). (In practice, not all of the unobserved variables have to be model parameters we are actually interested in; these can be integrated out).

Going beyond simple applications like the classifier mentioned above, handling the posterior gets difficult, though. Simply evaluating the posterior density $\theta \mapsto p(\theta | x)$ at single points is not enough in a Bayesian setting for usages such as prediction, parameter estimation, or evaluation of probabilities of continuous variables. The problem is that almost all of the relevant quantities depend on some sort of expectation over the posterior density, an integral of the form

$$\mathbb{E}[f(\Theta) | X = x] = \int f(\theta) p(\theta | x) d\mu(\theta), \quad (2.2)$$

for some measurable function f (with the base measure μ depending on the type of Θ). This in turn involves calculating the normalizing marginal

$$p(x) = \int p(x, \theta) d\mu(\theta). \quad (2.3)$$

in equation (2.1), often called the “evidence”.

When the distributions involved form a sufficiently “nice” combination, e.g., a conjugate pair (see Marin, Robert 2007, chapter 2.2.2; Murphy 2012, chapter 9.2.5), the integration can be performed analytically, since the posterior density has a closed form for a certain known distribution, or at least is a known integral. In general, however, this is not tractable, not even by standard numerical integration methods, and approximations have to be made. Even for discrete variables, the applicability of simple summation is limited by combinatorial explosion.

DIFFERENT TECHNIQUES for posterior approximation are available: among them are distribution-based approaches for general graphical models, such as variational inference (Murphy 2012, chapter 21 and 22) and other methods generalized under the

¹Note the abuse of notation regarding $p(\cdot)$; see page xi on notation.

framework of message passing (Minka 2005). The methods described in this thesis, however, fall into the category of Monte Carlo methods, and are based on sampling (Murphy 2012, chapter 23; Vihola 2020). Their fundamental idea is to derive, for a specified density of $\Theta \sim \pi$, a sampling procedure with a consistent estimator for expectations:

$$I^{(k)}(f) \rightarrow \mathbb{E}[f(\Theta)] = \int f(\theta)\pi(\theta) d\mu(\theta), \quad \text{as } k \rightarrow \infty \quad (2.4)$$

in some appropriate stochastic convergence (usually convergence in probability is enough). We leave out the conditional dependency on X in the following for simplicity of notation, and since the data are usually fixed in inference problems.

Examples of such methods are rejection sampling, importance sampling, and particle filters. Many Monte Carlo methods are defined in a form that directly samples a sequence of individual random variables $(Y^{(k)})_{k \geq 1}$, called a *chain*, for which the estimator is given by the arithmetic mean, such that a law of large numbers (LLN) holds:

$$I^{(k)}(f) = \frac{1}{k} \sum_{i=1}^k f(Y^{(i)}) \rightarrow \mathbb{E}[f(\Theta)] \quad (2.5)$$

If we can sample $Y^{(k)} \sim \pi$ exactly, they are i.i.d. and the LLN holds trivially; such samplers exist, but might also be difficult to derive or not possess good enough convergence properties (especially in high dimensions). Another large class of samplers is formed by *Markov Chain Monte Carlo* (MCMC) methods, which, instead of sampling exactly from the density, define $Y^{(k)}$ via a (time-homogeneous) Markov chain:

$$\begin{aligned} \mathbb{P}[Y^{(k+1)} \in dy \mid Y^{(k)} = y^{(k)}, \dots, Y^{(1)} = y^{(1)}] \\ = \mathbb{P}[Y^{(k+1)} \in dy \mid Y^{(k)} = y^{(k)}] \\ = K(dy \mid y^{(k)}) \end{aligned} \quad (2.6)$$

for all $k \geq 1$. By constructing the parametrized measure K , the *transition kernel*, in the right way, the resulting chain is ergodic with the target density π as the unique stationary distribution, i.e., for all measurable sets A ,

$$\int \pi(\theta) K(A \mid \theta) d\mu(\theta) = \int_A \pi(\theta) d\mu(\theta) = \mathbb{P}[\Theta \in A], \quad (2.7)$$

and the LLN for Markov chains holds. (For discrete spaces, this relation is more familiarly written as a left eigenvalue equation on a stochastic matrix: $\pi K = \pi$.) The advantage of MCMC methods is that they apply equally well to many structurally complex models, and treat densities in a uniform way, without requiring special knowledge about the specific distribution in question. I refer to Vihola (2020, chapter 6), Robert, Casella (1999), and Murphy (2012, chapters 24 and following) as introductions to MCMC theory and practice.

```

Start from an arbitrary  $Y^{(1)} = y^{(1)}$  with  $\pi(y^{(k)}) > 0$ 
for  $k \geq 1$  do
    Sample a proposal  $\hat{Y}^{(k)} \sim q(Y^{(k-1)}, \cdot)$ 
    With probability  $\alpha(\hat{Y}^{(k)}, Y^{(k-1)})$ , set  $Y^{(k)} = \hat{Y}^{(k)}$ ; else, keep  $Y^{(k)} = Y^{(k-1)}$ 
end for

```

Algorithm 2.1: General scheme for the Metropolis-Hastings algorithm.

FREQUENTLY, MCMC METHODS are variations of the *Metropolis-Hastings algorithm* (MH), which splits the general definition of the transition kernel into two parts: a proposal distribution, given by a conditional density q that needs to be easy to sample from, and an acceptance rate α . Subsequent samples are then produced by proposing values from q given the previous element of the chain, and incorporating them into the chain with a probability given through α (see algorithm 2.1). There exist many MH-based schemes with different properties and requirements: from the classical random-walk Metropolis algorithm with Gaussian proposals, over Reversible Jump MCMC for varying dimensions (Green 1995), to gradient-informed methods like Metropolis Adjusted Langevin and Hamiltonian Monte Carlo (HMC) (Betancourt 2018; Girolami, Calderhead 2011).

For multi-component structures, of the form $\Theta = [\Theta_1, \dots, \Theta_N]$, a good proposal distribution can be hard to find, though. One way to break down the problem is to use a family of componentwise updates, given by conditional distributions q_i operating on only one component of Θ , with the others fixed:

$$\begin{aligned}\hat{Y}_{-i}^{(k)} &= Y_{-i}^{(k-1)} \\ \hat{Y}_i^{(k)} &\sim q_i(Y_i^{(k-1)}, \cdot \mid Y_{-i}^{(k-1)})\end{aligned}\tag{2.8}$$

The components can be scalar or multivariate blocks, and the kernel may itself be any valid transition kernel (Vihola 2020, chapter 6.6). This allows one to freely mix different MCMC methods suitable for each variable in a problem.

This so-called “within-Gibbs” sampler bears its name because it is a generalization of the classical *Gibbs sampling* algorithm (S. Geman, D. Geman 1984): often, the simplest available set of transition kernels is given by the conditional densities $\theta_i \mapsto p(\theta_i \mid \theta_{-i}, x)$. They can directly be used as component proposals for a within-Gibbs sampler, leading to a cancelling acceptance rate of $\alpha \equiv 1$. This approach has the advantage of being very algorithmic, which makes it rather easy to apply, even by hand, to many models, and simply to express algorithmically. Hence, the method is a popular starting point for general probabilistic programming systems, most prominently BUGS (Lunn, Spiegelhalter, et al. 2009; Lunn, Thomas, et al. 2000) and JAGS (Plummer 2003; Plummer 2017).

In many real-world models, the factorization structure is quite sparse and results in small Markov blankets. Algorithms to derive Gibbs samplers exploit this large independency between variables. In short, they “trim” the dependency graph of the model to the local Markov blankets of each target variable, and derive either

a full conditional from it, where possible (for discrete or conjugate variables), or otherwise approximate it through appropriate local sampling (e.g., slice sampling) (see Plummer 2003).

As an example, consider a simple Gaussian mixture model with equal weights, specified as follows:

$$\begin{aligned}\mu_k &\stackrel{\text{iid}}{\sim} \text{Normal}(m, s) \quad \text{for } 1 \leq k \leq K, \\ Z_n &\stackrel{\text{iid}}{\sim} \text{Categorical}(K) \quad \text{for } 1 \leq n \leq N, \\ X_n &\stackrel{\text{iid}}{\sim} \text{Normal}(\mu_{Z_n}, \sigma) \quad \text{for } 1 \leq n \leq N.\end{aligned}\tag{2.9}$$

To derive the conditional distribution of Z_n given the remaining variables, we start by writing down the factorization of the joint density:

$$p(z_{1:N}, \mu_{1:K}, x_{1:N}) = \prod_k p(\mu_k) \prod_n p(z_n) \prod_n p(x_n | \mu_{z_n}).\tag{2.10}$$

From this, we can derive an unnormalized density proportional to the conditional by removing all factors not including the target variable:

$$p(z_n | z_{-n}, \mu_{1:K}, x_{1:N}) \propto p(z_n) p(x_n | \mu_{z_n})\tag{2.11}$$

This is equivalent to finding the Markov blanket of Z_n : only those conditionals relating the target variable to its children and parents remain. Since the clusters are drawn from a categorical distribution, the support is discrete, and we can find the normalization constant by summation:

$$\begin{aligned}p(z_n | z_{-n}, \mu_{1:K}, x_{1:N}) \\ = \frac{\text{Categorical}(z_n | K) \text{Normal}(x_n | \mu_{z_n}, \sigma)}{\sum_{k \in \text{supp}(Z_n)} \text{Categorical}(k | K) \text{Normal}(x_n | \mu_k, \sigma)},\end{aligned}\tag{2.12}$$

which can be expressed as a general discrete distribution over $\text{supp}(Z_n) = \{1, \dots, K\}$, with the unnormalized weights given by the numerator. Next, the conditionals of the μ_k have the form

$$\begin{aligned}p(\mu_k | z_{1:N}, \mu_{-k}, x_{1:N}) \\ \propto p(\mu_k) \prod_n p(x_n | \mu_k)^{\mathbb{1}(z_n=k)} \\ = \prod_n (\text{Normal}(\mu_k | m, s) \text{Normal}(x_n | \mu_k, \sigma))^{\mathbb{1}(z_n=k)}\end{aligned}\tag{2.13}$$

which we recognize as a product of conjugate pairs of normal distributions. More examples are extensively covered in Murphy (2012, chapter 24.2).

```

@model function normal_mixture(x, K, m, s,  $\sigma$ )
    N = length(x)

     $\mu$  = Vector{Float64}(undef, K)
    for k = 1:K
         $\mu[k]$  ~ Normal(m, s)
    end

    z = Vector{Int}(undef, N)
    for n = 1:N
        z[n] ~ Categorical(K)
    end

    for n = 1:N
        x[n] ~ Normal( $\mu[z[n]]$ ,  $\sigma$ )
    end

    return x
end

```

Listing 2.1: Turing.jl implementation of a Gaussian mixture model with prior on the cluster centers, equal cluster weights, and all other parameters fixed.

2.2 PROBABILISTIC PROGRAMMING

Probabilistic programming is a structured way implementing generative models, as described in the previous section, through the syntax of a programming language. It is beneficial to consider probabilistic programs not only as syntactic sugar for denoting the implementation of a joint probability density over some set of variables, but as organized objects in their own right: they open up possibilities that “black box” density functions cannot automatically provide. In more concise terms of J.-W. van de Meent et al. (2018):

Probabilistic programming is largely about designing languages, interpreters, and compilers that translate inference problems denoted in programming language syntax into formal mathematical objects that allow and accommodate generic probabilistic inference, particularly Bayesian inference and conditioning.

A probabilistic program differs from a regular program (that may also contain stochastic parts) through the possibility of being conditioned on: some of the internal variables can be fixed to observed values, from outside. As such, the program denotes on the one hand a joint distribution, that can be *forward sampled* from by simply running the program top to bottom and producing (pseudo-) random values. But at the same time, it also represents a conditional distribution, in form on the unnormalized conditional density, which together with an inference algorithm can also be *backward sampled* from. (Other terms, such as “evaluation” and “querying”, are used as well.) Consider the model (2.9) from above: to perform inference on it in Turing.jl (Ge, Xu, Ghahramani 2018), the probabilistic programming language

used in this thesis, its mathematical description might be translated into the Julia program given in listing 2.1.

We can then sample from the model in several ways using Julia:

```
julia> m = normal_mixture(x_observations, K, m, s,  $\sigma$ );
julia> forward = sample(m, Prior(), 10);
julia> chain = sample(m, MH(), 1000);
```

The value of `forward` will be an dataframe-like object containing 10 values for each variable sampled from the forward (i.e., joint) distribution, matching the size of `x_observations`. Similarly, `chain` will contain a length 1000 sample from a Markov chain targetting the posterior, conditionally on `x_observations`, created using the MH algorithm. If we were to write out code for these two functionalities manually, in idiomatic Julia, we would end up with at least two separate functions needed for the sampler:

```
function normal_mixture_sampler(N, K, m, s,  $\sigma$ )
     $\mu$  = rand(Normal(m, s), K)
    z = rand(Categorical(K), N)
    x = rand.(Normal.( $\mu$ [z], s))
    return  $\mu$ , z, x
end

function normal_mixture_logpdf( $\mu$ , z, x, K, m, s,  $\sigma$ )
    N = length(x)
     $\ell$  = 0.0
     $\ell$  += sum(logpdf(Normal(m, s),  $\mu$ [k]) for k = 1:K)
     $\ell$  += sum(logpdf(Categorical(K), z[n]) for n = 1:N)
     $\ell$  += sum(logpdf(Normal( $\mu$ [z[n]], s), x[n]) for n = 1:N)
    return  $\ell$ 
end
```

And still, with these, we would lack much of the flexibility that models written in Turing.jl: no general interface for sampling algorithms to automatically detect all latent and observed variables; no possibility for other, nonstandard execution forms as are needed for Variational Inference or gradient computation for HMC; no automatic name extraction and dataframe building for chains. All these points highlight the advantages of dedicated probabilistic programming languages (PPLs) over hand-written model code. (Additionally, there is of course a benefit of reducing errors introduced by the sampling function not matching the likelihood function, or errors involving log-probabilities.)

MANY PPLs ARE IMPLEMENTED as external domain-specific languages (DSLs), like Stan (Carpenter et al. 2017), JAGS (Plummer 2003), and BUGS (Lunn, Spiegelhalter, et al. 2009; Lunn, Thomas, et al. 2000). Others are specified in the “meta-syntax” of Lisp S-expressions, as Church (Goodman, Mansinghka, et al. 2012), Anglican (Wood, J. W. van de Meent, Mansinghka 2015), or Venture (Mansinghka, Selsam, Perov 2014). A third group is embedded into host programming languages with sufficient syntactic flexibility, for example Gen (Cusumano-Towner 2020) and Soss (Scherrer 2019) in Julia (besides the already named Turing.jl), or Pyro (Bingham et al. 2018) and PyMC3 (Salvatier, Wiecki, Fonnesbeck 2016) in Python.

The latter approach is advantageous when one wants to enable the use of regular, general-purpose programming constructs or interact with other functionalities of the host language. There are also a variety of further reasons why one would rather describe an inference problem in terms of a program than in more “mathematical” form, like as a graph or likelihood function. In a good probabilistic programming DSL, models will read as close to textbook model specifications as possible, while allowing to use the host language to:

- define recursive relationships,
- write models using imperative constructs, such as loops, or mutable intermediate computations for efficiency,
- optimize details of the execution, e.g. for memoization, likelihood scaling, or preliminary termination,
- use distributions over complex custom data structures, e.g. trees,
- perform inference involving complex transformations from other domains, for which implementations already exist, e.g. neural networks or differential equation solvers, or
- integrate calls to very complex external systems, e.g. simulators or renderers.

Depending on the choice of features should be supported, several possibilities for the implementation of such a DSL exist. All are based on some form of abstract interpretation. A rough distinction can be made between *compilation-based methods*, which statically translate the model code to a graph or density function, and *evaluation-based methods*, which dynamically or implicitly build such a structure at runtime, by allowing an inference algorithm to interleave the execution. The latter make it easier to include host-language control constructs. See J.-W. van de Meent et al. (2018) for a general introduction into some common implementation approaches for PPLs, and Goodman, Stuhlmüller (2014) for a detailed overview of the internals of one specific, continuation-based implementation called WebPPL (using a Lisp-based syntax).

`Models in Turing.jl` are written in `DynamicPPL.jl` syntax (Tarek et al. 2020), which transforms valid Julia function definitions into a reusable representation (`@model` is a Julia macro; see section 2.3 for more explanation). The result is a new function which produces instances of a structure of type `Model`, which in turn will contain the provided data, some metadata, and a nested function with the slightly changed original model code. In the concrete case of the model in listing 2.1, the resulting code would be approximately equal to the code in listing 2.2. The purpose of this is the following: the outer function, the “generator”, constructs an instance of the model for given parameters – usually done once per inference problem, to fix the observations and hyperparameters. Subsequently, the `sample` function can be applied to this instance with different values for the sampling algorithm, which in turn will use the evaluator function of the instance to run the model with chosen

```

function normal_mixture(x, K, m, s,  $\sigma$ )
  function evaluator((rng, model, varinfo, sampler, context, x, K, m, s,  $\sigma$ )
    N = length(x)
     $\mu$  = Vector{Float64}(undef, K)
    for k = 1:K
      dist_mu = Normal(m, s)
      vn_mu = @varname  $\mu$ [k]
      inds_mu = ((k,)),)
       $\mu$ [k] = tilde_assume(
        rng, context, sampler, dist_mu, vn_mu, inds_mu, varinfo
      )
    end
    z = Vector{Int}(undef, N)
    for n = 1:N
      dist_z = Categorical(K)
      vn_z = @varname z[n]
      inds_z = ((n,)),)
      z[n] = tilde_assume(
        rng, context, sampler, dist_z, vn_z, inds_z, varinfo
      )
    end
    for n = 1:N
      dist_x = Normal( $\mu$ [z[n]],  $\sigma$ )
      vn_x = @varname(x[n])
      inds_x = ((n,)),)
      if isassumption(model, x, vn_x)
        x[n] = tilde_assume(
          rng, context, sampler, dist_x, vn_x, inds_x, varinfo
        )
      else
        tilde_observe(
          context, sampler, dist_x, x[n], vn_x, inds_x, varinfo
        )
      end
    end
    return x
  end
  return Model(
    :normal_mixture, evaluator,
    (x = x, K = K, m = m, s = s,  $\sigma$  =  $\sigma$ ),
    NamedTuple()
  )
end

```

Listing 2.2: Slightly simplified macro-expanded code of the model in listing 2.1. The inner code is put into an `evaluator` closure, and every `tilde` statement is replaced by a `tilde_*` function, to which additional data and state information are passed.

sampler and context arguments, that are passed to the “tilde functions”, to which the statements of the form `expr ~ D` are converted.

A special distinction is made for the tilde functions of variables that are based on the model’s arguments. `DynamicPPL.jl` distinguishes between *assumptions*, i.e., latent variables that should be recovered through posterior inference, and *observations*, that need to be provided when instantiating the model and are conditioned upon. The latter by default will only contribute to the likelihood, instead of being sampled. But in certain cases, such as in probability evaluation or when using the complete model in a generative way, this behaviour can be different. For this purpose, the tilde functions for the variables `x[i]` in listing 2.2 are differentiated in a conditional statement.

Inside the tilde functions, the real stochastic work happens. Depending on the sampler and the context, values may be generated and stored in the `varinfo` object, and the joint log-likelihood incremented, as happens for most MCMC samplers. In this case, one call to the evaluator corresponds to one sampling step. In other situations, model evaluation serves the purpose of density evaluation, in which no new values need to be produced; this use case is needed for probability queries, or density-based algorithms (which might additionally use automatic differentiation on the density evaluation procedure). All shared information for external usage is thereby conventionally stored in the `varinfo` object, which resembles a dictionary from variable names² to values (internal sampler state can also be stored in the sampler object). Through the `sample` interface, the resulting values are then stored in a `Chains` object, a data frame containing a value for each variable at each sampling step.

From the point of view of a sampling algorithm, all that it sees is a sequence of tilde statements, consisting of a value, a variable name, and a distribution. `Turing.jl`, crucially, does not have a representation of model structure. This is sufficient for many kinds of inference algorithms that it already implements – Metropolis-Hastings, several particle methods, HMC and NUTS, and within-Gibbs combinations of these – but does not allow more intelligent usage of the available information. For example, to use a true, conditional, Gibbs sampler, the user has to calculate the conditionals themselves. Structure-based optimizations such as partial specialization of a model to save calculations, automatic conjugacy detection (Hoffman, Johnson, Tran 2018), or model transformations such as Rao-Blackwellization (Murray et al. 2017) cannot be performed in this representation.

2.3 COMPILATION AND METAPROGRAMMING IN JULIA

Julia (Bezanson, Edelman, et al. 2017) is a programming language with a strong, dynamic type system with nominal, parametric subtyping and elaborate multiple dispatch. It uses LLVM (LLVM Project 2019) for JIT-compilation and while it is

²These `VarName` objects, constructed by the macro `@varname`, simply represent an indexed variable through a symbol and a tuple of integers.

dynamically typed, a combination of method specialization and type inference allows it to produce very optimized, fast machine code (Bezanson, Chen, et al. 2018). The language is syntactically designed to bear a certain resemblance to Matlab, Python, or Ruby, but contrary to them, it is its own compiler, and not primarily the reliance on libraries calling foreign functions (e.g., Numpy), which is intrinsically enabling C-like speed. Although Julia does rely on, e.g., BLAS and LAPACK for numerical algebra, there is nothing that fundamentally prevents implementing their functions: true array types, fast loops, and various optimiations are available, as opposed to languages like Python, which are fundamentally limited by to their dynamic interpretation. This advantage carries over to domains outside of numeric computation, of course.

On top of that, the language is built on a very open compilation model. Underlying the surface syntax is an abstract syntax tree (AST), that is used internally to the compiler, but also exposed to the programmer through macros, which allow to transform pieces of code at compile time. These macros resemble proper hygienic, LISP-style code transformations (cf. Hoyte 2008), not simple text-substitutions as C preprocessor macros. As an example, look at the following method³ that sums up the sin values of a list of numbers:

```
function foo(x)
    y = zero(eltype(x))
    for i in eachindex(x)
        @show y += sin(x[i])
    end
    return y
end
```

The invocation of the standard library macro `@show` will be treated by the compiler, during parsing, as a function call receiving as input the following data structure, representing `y += sin(x[i])` in S-expression-like form:

```
Expr(:(+=), :y, Expr(:call, :sin, Expr(:ref, :x, :i)))
```

In this particular case, the nested structure is not taken advantage of or transformed, but simply converted to a string used to print the value of the expression, labelled by its form in the code:

```
macro show(ex)
    blk = Expr(:block)
    unquoted = sprint(Base.show_unquoted, ex) * " = "
    assignment = Expr(:call, :repr, Expr(:(=), :value, esc(ex)))
    push!(blk.args, Expr(:call, :println, unquoted, assignment))
    push!(blk.args, :value)
    return blk
end
```

³The terminology of Julia uses *function* for a callable object, which can have multiple *methods* for different combinations of argument types. This is what allows multiple dispatch: when a function is applied, the types of the arguments are determined, and the most specific matching methods selected and called. For example, the `+` function has many methods for adding integers, floats, arrays, etc.

```

1: (%1::Core.Compiler.Const(foo, false), %2::Array{Float64,1})
   %3 = eltype(%2)::Compiler.Const(Float64, false)
   %4 = zero(%3)::Float64
   %5 = eachindex(%2)::Base.OneTo{Int64}
   %6 = iterate(%5)::Union{Nothing, Tuple{Int64,Int64}}
   %7 = (%6 == nothing)::Bool
   %8 = not_int(%7)::Bool
   br $3 (%4) unless %8
   br $2 (%6, %4)
2: (%9, %10)
   %11 = getfield(%9, 1)::Int64
   %12 = getfield(%9, 2)::Int64
   %13 = getindex(%2, %11)::Float64
   %14 = sin(%13)::Float64
   %15 = (%10 + %14)::Float64
   %16 = repr(%15)::String
   %17 = println("y += sin(x[i]) = ", %16)
   %18 = iterate(%5, %12)::Union{Nothing, Tuple{Int64,Int64}}
   %19 = (%18 == nothing)::Bool
   %20 = not_int(%19)::Bool
   br $3 (%15) unless %20
   br $2 (%18, %15)
3: (%21)
   return %21

```

Listing 2.3: SSA-form of the lowered form of the method `foo(::Vector{Int})` as defined defined above, annotated with inferred types (as through `@code_warntype`).

The result is then spliced back into the AST, which is compiled further as if it were written as

```

function foo(x)
    y = zero(eltype(x))
    for i = eachindex(x)
        begin
            println("y += sin(x[i]) = ", repr(var"#1#value" = (y += sin(x[i]))))
            var"#1#value"
        end
    end
    return y
end

```

(Note the automatic conversion of the symbol `:value` to a generated name `#1#value`, in order to not possibly shadow any variables from the calling scope.)

After macro expansion, the code of the method is *lowered* into an intermediate representation consisting of only function calls and branches. This comprises of several transformations: for one, certain syntactic constructs are “desugared” into primitive function calls. For example, array accesses, `x[i]`, are replaced by calls to the library function `getindex(x, i)`. The for loop in the example is converted into a while loop using the `iterate` library function:

```

iterable = eachindex(x)
iter_result = iterate(iterable)
while !(iter_result == nothing)
    i, state = iter
    @show y += sin(x[i])
end

```



```

    iter_result = iterate(iterable, state)
end

```

Consequently, all nested expressions are split apart, so that only simple, unnested calls remain, and any subsequent assignments to variables are linearized to a series of definitions, with newly introduced names of the form %i. The remaining control flow statements (e.g., while loops and conditionals) are represented through sequence of labelled *basic blocks*, with (possibly conditional) jumps between them. The sequence of assignments is further processed into *single static assignment (SSA) form* (Singer 2018), the characteristic property of which is that every variable is assigned exactly once, thus giving it a unique, position-independent name to each intermediate value. By introducing this immutability guarantee, the resulting code is, in a certain sense, referentially transparent, which facilitates data-flow analysis, and makes many transformations easier. Accordingly, SSA form is widely used in intermediate forms of compiler systems, simplifying transformations and optimizations. The result of the translation of our example into three basic blocks can be found in listing 2.3.

There is one notable complication regarding conversion to SSA form: we need to be able to distinguish between assignments of variables arising from “joined” control flow. Consider the assignment of *y* in the following code example:

```

x = f()
if !g()
    y = x - 1
else
    y = x + 1
end
h(y)

```

Here, the value of *h(y)* depends on two possible locations of *y* – hence, we cannot simply rename every variable in a naive way. Instead, in the variant of SSA form used in this text and most of Julia, values of variables that are assigned in multiple parent blocks are passed on as *block arguments*, as in figure 2.1 on the right, and subsequently in this work. This makes basic blocks resemble local function, in a way, and cleanly resolves the problem of joins just like functions handle variable inputs. The traditional, functionally equivalent alternative is to introduce ϕ -functions (Rosen, Wegman, Zadeck 1988), which are defined ad-hoc to distinguish between several values depending on the control path taken before. This form is shown in the same figure on the left.

Note that until now, the operations involved were purely syntactic in nature, and could be performed by solely taking into account the code of the function *foo*. As soon as *foo* is called on a concrete type during evaluation, though, the most specific method fitting to the argument types will be selected, and type inference on its body be applied. If we go on and call *foo*([1, 0]), with *Vector{Int}* as the sole argument type, the types as annotated in the same listing will be inferred.

The last step of compilation within Julia consists of inlining and optimizing the typed intermediate code, resulting in the form shown in listing 2.4. There, several called methods have been inlined, and concrete argument types to invoke been

```

1 -- %1 = arraysize(x, 1)::Int64
   |   %2 = slt_int(%1, 0)::Bool
   |   %3 = ifelse(%2, 0, %1)::Int64
   |   %4 = slt_int(%3, 1)::Bool
   |___ goto $3 if not %4
2 --   goto $4
3 --   goto $4
4 -- %8 = φ ($2 => true, $3 => false)::Bool
   |   %9 = φ ($3 => 1)::Int64
   |   %10 = φ ($3 => 1)::Int64
   |   %11 = not_int(%8)::Bool
   |___ goto $22 if not %11
5 -- %13 = φ ($4 => 0.0, $21 => %18)::Float64
   |   %14 = φ ($4 => %9, $21 => %42)::Int64
   |   %15 = φ ($4 => %10, $21 => %43)::Int64
   |   %16 = arrayref(true, x, %14)::Float64
   |   %17 = invoke sin(%16::Float64)::Float64
   |   %18 = add_float(%13, %17)::Float64
   |   %19 = sle_int(1, 1)::Bool
   |___ goto $7 if not %19
6 -- %21 = sle_int(1, 0)::Bool
   |___ goto $8
7 --   nothing::Nothing
8 -- %24 = φ ($6 => %21, $7 => false)::Bool
   |___ goto $10 if not %24
9 --   invoke getIndex()::Tuple, 1::Int64)::Union{}
   |___ $(Expr(:unreachable))::Union{}
10 -   goto $11
11 -   goto $12
12 -   goto $13
13 -   goto $14
14 - %32 = invoke :(var"#sprint#339")(
   |       nothing::Nothing, 0::Int64, sprint::typeof(sprint),
   |       show::Function, %18::Float64
   |___ )::String
   |       goto $15
15 -   goto $16
16 -   goto $17
17 -   invoke println("y += sin(x[i]) = "::


---



```

Listing 2.4: Typed and optimized code of the call `foo([1.0])` in SSA form, as obtained through `@code_typed` (the extra bars are due to the formatting of `CodeInfo`).

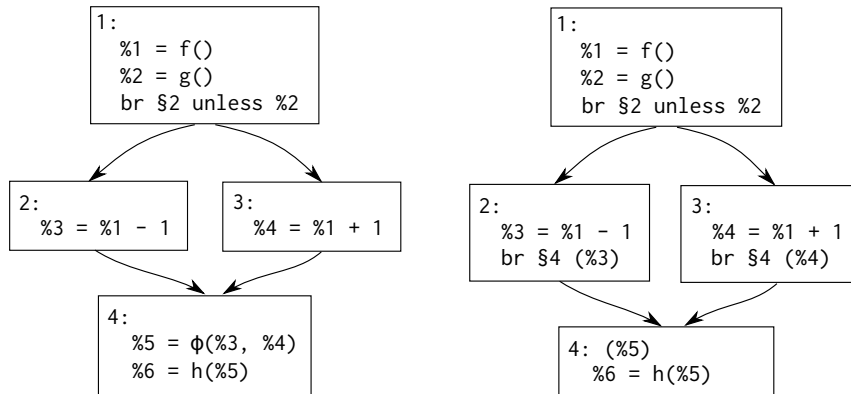


Figure 2.1: Two control flow graphs of the same function, illustrating the correspondence between SSA representations using ϕ -functions and block arguments. The SSA variables %3 and %4 correspond to the values of y in the two branches, which are merged in %5.

inferred. This is in true, traditional SSA form, with all variable slots eliminated, and block arguments converted to the mentioned ϕ -functions. Finally, this representation will be translated and sent to LLVM for compilation, where further optimization can happen, and machine code will be generated and executed, as well as stored for later usage as part of the just-in-time compilation mechanism.

A KEY PRINCIPLE in Julia’s compilation model is type specialization (Bezanson, Chen, et al. 2018). As we have seen, whenever a function call is evaluated, the types of the arguments are first determined, and then the most specific method selected and called. This automatically gives the language dynamic semantics: an implementation can perform evaluation at every call. In practice, however, at this point multiple dispatch and JIT compilation combine into one of the main principles of optimization. Instead of evaluating the same code over and over again, methods are JIT-compiled the first time they are called. The compiled code is then cached in the method table. Method compilation does not happen recursively at once, though. Only when the body of a compiled method is then executed with concrete arguments, the same process is performed again, for each invoked method.

So, in a sense, JIT compilation can be seen as a function that returns compiled code, given a function and a tuple of types. Similar to macros, which transform original code, given an expression, this process of generating compiled methods from types is customizable in Julia: via so-called *generated functions*, there exists a process to dynamically generate code, given argument types – a form of stated programming (Bolewski 2015; Rompf, Odersky 2010). Such generated functions, when called, are not directly translated into machine code: instead, they emit new code to the compiler, based on the types of their arguments. The new code is then JIT-compiled. For example, when we have two methods of a function f :

```
f(x::Int) = println("Int")
f(x::String) = println("String")
```

we could replace them with the following generated function:

```
@generated function f_generated(x)
  if x == Int
    return :(println("Int"))
  elseif x == String
    return :(println("String"))
  else
    error("Method error")
  end
end
```

Calling `f_generated(1)` will then determine the argument type (`typeof(1) == Int`), and pass it to the function body of `f_generated`. There, the conditional will select the first branch, and the expression `:(println("Int"))` be returned. This is now passed back to the compiler, which will lower the code and compile the method for `Int` arguments, and store the result in the method table. The stored code can then be executed – on the arguments that were used to determine the type tuple the generated function has been called with! The next time `f_generated` is executed, the function body is *not* executed anymore, but the generated code of `:(println("Int"))` directly looked up⁴. Of course, simply replacing dispatch, as with this example, is not what generated functions are used for in practise. Most applications concern parametric types with statically known shape arguments, such as tuples, named tuples, or array ranks. They can also be used for type-level computations on values that become known only known at runtime, through singleton types such as `Val`.

The direct generation of code, given argument types, is however not the furthest we can go. For onegenerated functions are not only allowed to return `Expr` objects – the internal representation of the surface AST – but also `CodeInfo` objects, which are the internal representation of lowered code in (almost) SSA form. This, on its own, would not be of much use most times, but there is a second, more interesting feature: it is possible to query the `CodeInfo` of a method by reflection, given a function and an argument type tuple. Combining these two, we now have all the tools to implement IR-level code transformations as follows:

1. Define a generated function, taking as arguments another function and its arguments.
2. Within this function, obtain the IR of the method of the passed-in function for the remaining arguments.
3. Transform this IR however necessary.
4. Return the IR, which will now be compiled and called on the actual arguments.

Importantly, unlike macros, such transformation can be performed *recursively*: by inserting the same generated function to the inner function calls in step 3. Since the transformation operates not during parsing, the function to be transformed

⁴A caveat: technically, the compiler is still free to call the generating code multiple times – which is the reason generated functions should never involve side effects or depend on external state.

needs not be known beforehand, and not be present literally in the code – the generated function can be called on every available callable object, at any time during runtime. This makes it possible to transform even functions from other libraries, internally calling yet other functions. One particular example of this principle is source-to-source automatic differentiation, as shown in the next chapter: a call to a function `gradient(f, x, y)` will obtain the IR of the method for `f` on `typeof(x)` and `typeof(y)`, produce differentiated code, and call the result on `x` and `y`. Naturally, differentiating `f` involves recursively differentiating the other, unknown functions within it, too (down to “primitive” functions, whose derivative is known), and combining the results using the chain rule.

This metaprogramming pattern is extremely powerful, and becoming more and more popular. It allows to change evaluation semantics in more profound ways than multiple dispatch can: by rewriting the code of the called function, it is possible to change what invoking a method within its body means. Through this, several abstract interpretation algorithms can be realized, by extending the existing data path with additional metadata (such as AD, or concolic execution), or non-standard execution be implemented (e.g., continuation-passing style transformations). There exist already two Julia packages with the goal of simplify working with this kind of transformation: `Cassette.jl`⁵, which provides overloadable function application by a so-called “overdubbing” mechanism, abstracting out some common patterns; and `IRTools.jl`⁶, which has a more user-friendly alternative to `CodeInfo`, and a macro similar to `@generated` that makes writing recursive functiona IR-transformation using this data structure easier. The latter is what the work of this thesis builds on.

2.4 AUTOMATIC DIFFERENTIATION AND COMPUTATION GRAPHS

This section may appear as an outlier from the original topic; in fact, the mathematics developed here are not even used later. However, it is important to explain the interrelations between automatic differentiation, computation graphs, and IR transformations, to be able to understand how SSA-form representation is a natural structure for extracting and analyzing computation graphs, and how the necessary transformations arise in practise. To appreciate how the form of the computation graphs interacts with the mathematics, some foundations need to be introduced first.

Many algorithms in machine learning and other domains can be expressed as an optimization problem over a multivariate function with scalar output – typically a loss function over a parameter space, which measures how far a model prediction is from the true target values. The optimal model is then just that one for which the parameters minimize the loss function. When the loss function is (sub-)differentiable, there exist a variety of gradient-based optimization methods to find this optimum (or, in the non-convex case, at least a practically sufficient local minimum).

⁵<https://github.com/jrevels/Cassette.jl>

⁶<https://github.com/FluxML/IRTools.jl>

While in some cases the loss function is simple enough to find the gradient by hand, in general, the model, and therefore the loss function, may be specified in terms of rather complicated programs, for which hand-writing derivatives is difficult to infeasible. For this reason, computerized methods for differentiation have been developed (Baydin et al. 2018; Griewank, Walther 2008). These can be categorized into three classes:

- Finite differences
- Symbolic differentiation
- Automatic differentiation

In finite differences, the idea is to discretize the definition of derivatives, and numerically evaluate the function within an environment. This is simple to implement, but does not scale well with the dimension of the involved space, and can become numerically unstable in various ways (Press et al. 2007, section 5.7). Symbolic differentiation works through representing the functions in question as symbolic algebraic objects, and applying the differentiation rules as one would manually. This does not lose precision or introduce divergence, but can suffer from blow-up of the size of the generated expressions; additionally, it requires the functions to be expressed in a custom representation, different from normal functions or programs.

AUTOMATIC DIFFERENTIATION, the third category, is perhaps unfortunately named. The idea is to not start from functions as black-box or symbolic objects, but from programs. Then the perturbation that makes up the value of the derivative at a point is propagated through the steps of the program. For this to work, there needs to exist an explicit representation of the computation graph at ov the evaluation at a point, which is what makes the topic relevant for this thesis.

To understand how this works, let us first start with the mathematics. What is a derivative, really? When we talk about gradients, which is what we really need in a gradient algorithm, this is usually a rather informal term for “the vector of partial derivatives”, which then points into an ascent direction. This is however not the most natural form to work with in a compositional approach. Instead of starting with a limit of tangents, more insight is provided by viewing derivatives as best-approximating linear operators. One of the most general definitions is provided through the *Fréchet derivative* (Bronstein, Semendjajew 1995, p. 463), essentially a generalization of the total differential. Let X and Y be normed spaces. A function $f : U \subseteq X \rightarrow Y$ is Fréchet differentiable at a point $x \in U$ if there exists a bounded linear operator $A : X \rightarrow Y$ such that

$$\lim_{\|h\|_X \rightarrow 0} \frac{\|f(x+h) - f(x) - A(h)\|_Y}{\|h\|_X}. \quad (2.14)$$

When such an A exists, it is unique, and we may call it *the* derivative of f at x , writing $Df(x) = A$. When the derivative exists for all x , we can use D as a well-defined higher-order function on its own; we will assume this in the following.

The important fact here is that $Df(x)$ is still a function: specifically, a linear function approximating how f reacts to a perturbation, Δ , around x . Or, in other words:

$$f(x + h) = f(x) + Df(x)(h) + o(\|h\|). \quad (2.15)$$

This fact allows one to propagate differential values through composed functions, by the chain rule, which we write in the following compositional form:

$$D(\phi \circ \psi)(x) = D\phi(\psi(x)) \circ D\psi(x) \quad (2.16)$$

In the one-dimensional case, we simply have

$$D\phi(x) = \Delta \mapsto \partial_1 \phi(x) \Delta, \quad (2.17)$$

since linear maps are exactly multiplications by a scalar. Therefore, we can recover

$$\begin{aligned} D(\phi \circ \psi)(x)(\Delta) &= (D\phi(\psi(x)) \circ D\psi(x))(\Delta) \\ &= D\phi(\psi(x))(D\psi(x)(\Delta)) \\ &= D\phi(\psi(x))(\partial_1 \psi(x) \Delta) \\ &= \partial_1 \phi(\psi(x)) \partial_1 \psi(x) \Delta \\ &= (\partial_1 \phi(\psi(x)) \partial_1 \psi(x)) \Delta, \end{aligned} \quad (2.18)$$

as we know it from calculus. Here, the product in the resulting expression arises from the fact that we propagated through $\partial_1 \psi(x) \Delta$ as the input value of $D\phi(\psi(x))$. It is remarkable that this formula is not entirely compositional, though: to construct $D(\phi \circ \psi)$, it is not only necessary to know $D\phi$ and $D\psi$, but also ψ (Elliott 2018). This is not as bad as it seems, though: as we see below, automatic differentiation algorithms evaluate both $((\phi \circ \psi)(x))$ and $D(\phi \circ \psi)(x)$ at once, in lockstep fashion, so that the intermediate values of the former can be reused in calculation of the latter.

Consider the specific case of $f(x, y) = \sin(x) - y$. We will write $g = (x, y) \mapsto x - y$ instead of the infix subtraction operator, with a derivative of $Dg(x)(\Delta_1, \Delta_2) = \Delta_1 - \Delta_2$. We have:

$$\begin{aligned} Df(x, y) &= D(g \circ (\sin \otimes \text{id}))(x, y) \\ &= Dg((\sin \otimes \text{id})(x, y)) \circ D(\sin \otimes \text{id})(x, y) \\ &= ((\Delta_1, \Delta_2) \mapsto \Delta_1 - \Delta_2) \circ (D\sin(x) \otimes D\text{id}(y)) \\ &= (\Delta_1, \Delta_2) \mapsto \partial_1 \sin(x) \Delta_1 - 1\Delta_2 \\ &= (\Delta_1, \Delta_2) \mapsto \cos(x) \Delta_1 - \Delta_2. \end{aligned} \quad (2.19)$$

In order to do this programmatically,

TO CONSTRUCT the full gradient from this, however, requires to evaluate $Df(x)$ at all the N basis functions. In order to overcome this, we can rewrite the above equation as

$$\begin{aligned} \partial_i f(x) &= Df(x)(e^i) \\ &= \langle 1, Df(x)(e^i) \rangle \\ &= \langle D^* f(x)(1), e^i \rangle, \end{aligned} \quad (2.20)$$

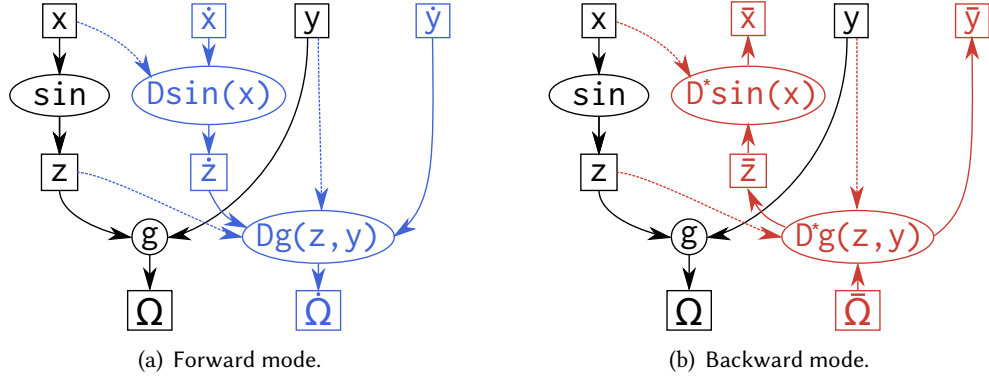


Figure 2.2: Computation graph and intermediate expressions of the expression $g(\sin(x), y)$, together with the derivative graphs in forward- and backward mode. Dashed arrows indicate re-use of primal values in the derivative graph.

where $D^*f(x)$ is the *adjoint operator*.

If $f : U \subseteq \mathbb{R}^N \rightarrow \mathbb{R}$, as is the case for loss functions of parametric models, the partial derivatives can be recovered individually by applying the Fréchet derivative to basis vectors e^i :

$$Df(x)(e^i) = \partial_i f(x) \quad (2.21)$$

This corresponds to the fact that the partial derivatives are the sensitivities to perturbation in coordinate directions – a special case of directional derivatives (which can be recovered the same way by application of the differential to any vector with unit norm.)

When we want to recover the gradient from the Jacobian, we can either multiply it with all the basis vectors, or simply transpose it:

$$\begin{aligned} Df(x)(e^i) &= \langle 1, [\cos(x), -1]^T e^i \rangle \\ &= \langle [\cos(x), -1], e^i \rangle \end{aligned} \quad (2.22)$$

(Note that the adjoint operator inverts the order of composition)

$$\begin{aligned} D^*f(x, y)(\delta) &= D^*(g \circ \sin \otimes \text{id})(x, y)(\delta) \\ &= (D^* \sin \otimes D^* \text{id} \circ D^* g)(x, y)(\delta) \\ &= ((D^* \sin(x) \otimes D^* \text{id}(y)) \circ (\xi \mapsto [\xi, -\xi]))(\delta) \\ &= ((\partial \sin(x) \times \cdot) \otimes (1 \times \cdot))([\delta, -\delta]) \\ &= [\cos(x), -1] \cdot \delta \end{aligned} \quad (2.23)$$

3 Implementation of Dynamic Graph Tracking in Julia

It has been mentioned above that there is a trade-off between source-transformation methods and library-based approaches for tracking computation graphs. Since the ultimate goal of this work was to analyze dynamic probabilistic models written in `Turing.jl`, properties of both were desired. Inspired by the work of Innes (2018), it seemed most promising to start from a source-transformation based approach implemented over the intermediate representation, especially from a usability point of view. The advantages of using IR over the surface AST are the same: there is less overhead from handling multiple syntactic forms, and naming is already referentially transparent. Additionally, there are existing Julia packages to simplify handling the IR data structures and set up the transformations.

However, the dynamicity of the trace structure of general probabilistic programs needs to be preserved and exposed to the user, for each function evaluation – which is different from the AD usage, where the adjoint function is already the ultimate goal, and does not change with the arguments. Hence, I developed a method for a hybrid version: through an IR transformation, the original code of a function to be tracked should be extended by additional statements to record a trace of the executed statements and control flow operations at runtime. The algorithm and data structure on which this approach is based have already been shortly described in Gabler et al. (2019), and will be more extensively explained below. An open source implementation is available online¹.

As we have seen above, in section 2.3, generated functions allow the inspection and transformation of the intermediate representation passed-in functions. This technique can be applied to recursively traverse the implementation of a given function, annotating each operation with necessary tracking statements, and changing the inputs and outputs accordingly to extract this information from outside. To ensure sufficient generality, we require the following properties of the tracking system:

1. Storage of all intermediate values during execution.
2. Symbolic capture expressions and branches in an analyzable, graphical form.
3. Preservation of the relation of each part of the structure to the corresponding

¹<https://github.com/TuringLang/IRTracker.jl>

original IR.

4. Proper nesting of this information for nested function calls, making relations between arguments and function inputs recoverable.
5. Correct handling of constants and primitive functions in the IR.
6. Extensibility of the tracking functions, to allow multiple possible ways to analyze code (e.g., by different definitions of what should be recorded).
7. A way to add custom metadata to the recorded structure during tracking.

This kind of operation will be similar to the (explicit) construction of Wengert lists in backwards-mode AD (see section 2.4); but contrary to there, the nested call structure and control flow shall be preserved as well. Hence, we call this structure *extended Wengert list*.

3.1 EXTENDED WENGERT LISTS

The extended Wengert list structure is implemented in Julia through nested objects of an abstract supertype `AbstractNode`, with several concrete subtypes for the different kinds of nodes. Additionally, there are special types for the tape- and block references, and an expression type `TapeExpression`, mimicking the built-in `Expr`, but adding more semantic distinctions (such as between references and constants, and between primitive and non-primitive function calls). On top of this, an API to query the graph structure is provided, allowing, for example, to find all children or parents of a tape reference up to a certain depth, or extract data from nodes, such as referenced variables, arguments, or metadata.

Figure 3.1 illustrates the resulting extended Wengert list for one run of a short stochastic function:

```
geom(n, beta) = rand() < beta ? n : geom(n + 1, beta)
```

(for readability, recorded to only three levels of nesting). The function draws a sample from the geometric distribution with parameter `beta`, starting to count at value `n`. On the left, we have its IR in textual form, consisting of two blocks. The central part is the graph of nested nodes. There, values and jumps from the top-level call are recorded in their encountered order, as nodes with “tape references” @1 to @9. SSA variables (%i) occurring in expressions of SSA definitions are also replaced in the nodes by the respective tape references. Each node is linked to the original IR statement it records, as indicated by the red arrows.

In the lower middle part, we see the node corresponding to the statement `%7 = geom(%6, %3)`. It is recorded at reference @8 with expression `geom(@7, @3)` and value 4 (the notation `<geom>(@7, @3, ())...` indicates that `geom` is a constant, and no variadic arguments are passed). The values of the arguments of this call can be inspected by looking up the respective references. Since `geom` is not a primitive function, the node holds tape of child nodes as well. In this case, it is equivalent to the top level, due to the recursivity of `geom`. We can see the three arguments @1, @2, and @3, corresponding to the block arguments %1, %2, and %3, with the value of @2

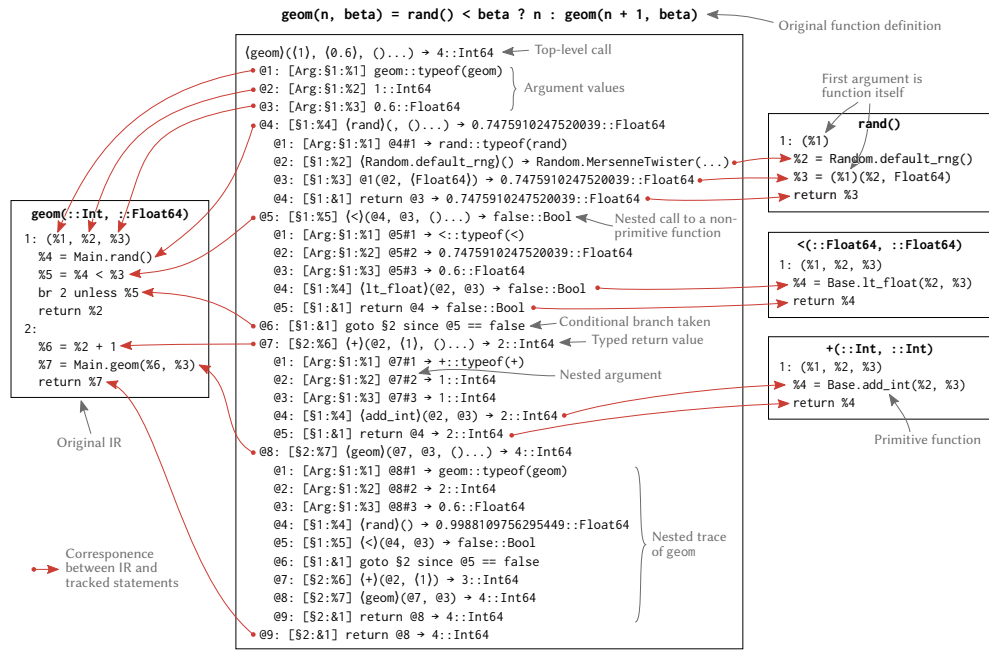


Figure 3.1: Extended Wengert list for one run of the stochastic function `geom` (only three levels shown). The central box is the tracked graph of the call `geom(1, 0.6)`. The other boxes show the original IR of the called non-primitive functions, to which the nodes are linked. Angle brackets indicate constant values.

being now 2 instead of 1. Further we can see function calls of `rand` and `<` as well as a conditional jump, corresponding to the branch the original IR, followed by calls of `+` and `geom`. Following back the tape references from the result value @9, the data path of the trace can be extracted. It can be used for reverse-mode AD, and only these nodes would be recorded in a conventional Wengert list. In our system, however, we also record the nodes on the control path, consisting of @6 and the nodes it depends on.

3.2 AUTOMATIC GRAPH TRACKING

Recording an extended Wengert list requires to record all block arguments, SSA definitions, and taken branches, with their actual values and metadata. This is achieved by extending the IR with new statements creating nodes and recording them on the extended Wengert list structure described above. Care needs to be taken to properly record function calls, since we need to ensure that non-primitive functions are recursively tracked.

The transformed code of the example function `geom`, whose IR is displayed in figure 3.1 above, is displayed in listing 3.1. First, a “graph recorder” object is set up in the extra argument %5. In this, the original IR is stored. Subsequently, every

statement is replaced by a call to one of the trackedX functions, to which both the function and its arguments, wrapped into TapeExpressions directly (for constants) or indirectly (through trackedvariable and trackedargument, which preserve the symbolic mapping to SSA variables). The record! function takes care of constructing the child node of the possibly nested call, and storing them on the recorder object.

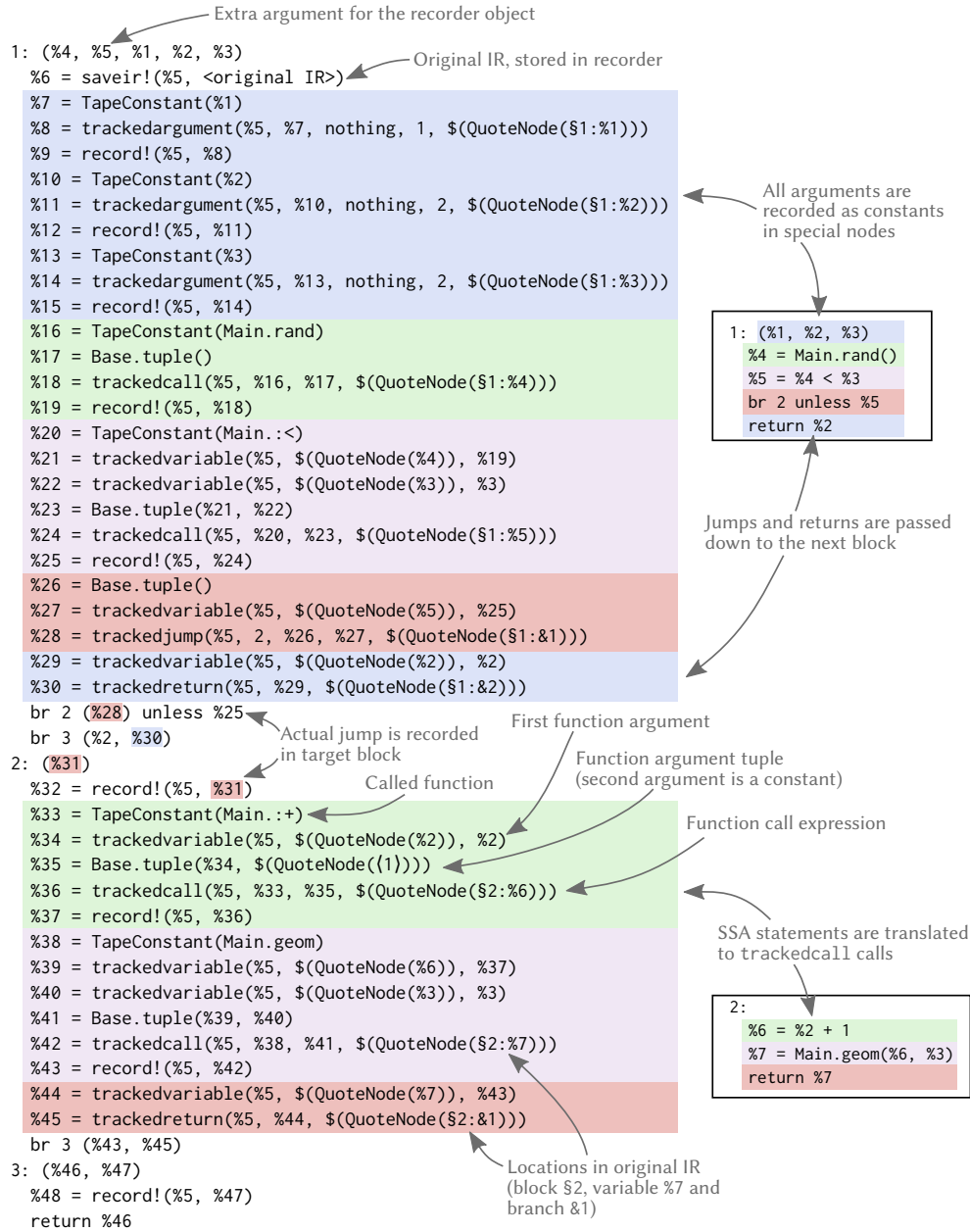
Branches, tracked with trackjump and trackedreturn, cannot be stored on the recorder object before they are taken, of course. The solution is to first construct the respective nodes of all possible branches of a block, and adding them as an extra argument to the branches taken. Then, in each target branch, the jump node from which the branch originated is recorded immediately. As a special case, all return branches are converted to unconditional jumps to one new block at the end, which contains a single unified return statement. This way, they can be treated in the same way as other branches.

The resulting IR consists of about three to five times as many statements as the original. The transformation, due to JIT compilation, is performed at most once per method and then stored as compiled code. However, the tracking – the recording of all statements in the extended Wengert list structure – happens at every execution during runtime.

Please see the appendix for a pseudo-code specification of the IR transformation in Algorithm 1, and the transformed code for the geom function in Figure 3.

make sure the figure is on same spread as explanation

describe algorithm



Listing 3.1: Tracked IR of the method `geom(::Int, ::Float64)`. Parts corresponding to original IR are highlighted in matching colors.

This transformation happens inside a generated function called by `trackcall`, which assembles the resulting value and IR into a new node with the correct metadata.

Missing from the description are the recording of metadata, the exact constructions of nodes, and the mechanisms to correctly rename SSA variables during the transformation and tape references at runtime.

```

function trackcall(ir)
  Initialize new_ir                                     ▶ create empty IR object
  for old_block in blocks (ir) do
    Add an empty block new_block to new_ir
    if (arg) is the first block then
      Add variable % recorder to new_block
    end if
    for arg in arguments ( old_block ) do
      Add arg to new_block
    end for
    if there exist branches to old_block then
      Add argument % branch_node to new_block
      Add statement to new_block , recording % branch_node in % recorder
    end if
    for arg in arguments ( old_block ) do
      Add statement % node to new_block , creating a node for arg
      Add statement to new_block , recording % branch_node in % recorder
    end for
    for stmt in statements ( old_block ) do
      if stmt is a normal call then
        Add statement % call_node to new_block , calling trackcall on stmt
        Add statement to new_block , recording % call_node in % recorder
      else if stmt is a “special” call or constant then
        Add statement % node to new_block , creating a node for stmt
        Add statement to new_block , recording % node in % recorder
      end if
    end for
    for branch in branches ( old_block ) do
      if branch is a return branch then
        ▶ Substitute return by a branch to the “return block”
        Add statement % return_node to new_block , creating a return node correspond-
ing to branch
        Add branch to new_block , targeting the return block, copying branch ’s argument,
with % return_node as extra argument
      else
        Add statement % branch_node to new_block , creating a node for branch
        Add branch to new_block , copying branch , with % branch_node as extra argu-
ment
      end if
    end for
  end for

  ▶ Set up “return block”
  Add block return_block to new_ir
  Add arguments % return_value , % return_node to return_block
  Add statement to return_block , recording % return_node in % recorder
  Add statement % result to return_block , creating a tuple of % return_value and
recorder
  Add return to return_block , returning % result
end function

```

3.3 EVALUATION

List of Algorithms

2.1	General scheme for the Metropolis-Hastings algorithm.	8
3.1	IR transformation to record an extended Wengert list (simplified) . . .	30