

ORACLE®

# Compilation of Probabilistic Programs

Jean-Baptiste Tristan

Principal Member of Technical Staff

Machine Learning Research Group - Oracle Labs

October 6, 2018

# Safe Harbor Statement

The following is intended to outline our general product direction. It is intended for information purposes only, and may not be incorporated into any contract. It is not a commitment to deliver any material, code, or functionality, and should not be relied upon in making purchasing decisions. The development, release, and timing of any features or functionality described for Oracle's products remains at the sole discretion of Oracle.

# This Talk

For machine learning people:

Compilers are marvels of engineering,  
we can borrow ideas and principles from 60 years of research

For compiler people:

Probabilistic programming compilation brings new  
interesting and challenging problems

# Compilers and Interpreters

An interpreter *simulates* a program

A compiler *translates* the program into machine code

# Why Care about Compilers?

We (hopefully) agree that probabilistic programming is useful

But what about *compiling* probabilistic programs?

Isn't an interpreter good enough?

In 2013, we designed and implemented the Augur language:

MAP inference for Bayesian networks *compiled* to GPU code

Running LDA:

- Augur (compiled): 1 GB, less than 20 minutes on GPU
- FACTORIE (interpreted): More than 6 hours on CPU
- Jags (interpreted): More than 128 GB, failed

# Why Care about Compilers?

We (hopefully) agree that probabilistic programming is useful

But what about *compiling* probabilistic programs?

Isn't an interpreter good enough?

In 2013, we designed and implemented the Augur language:

MAP inference for Bayesian networks *compiled* to GPU code

Running LDA:

- Augur (compiled): 1 GB, less than 20 minutes on GPU
- FACTORIE (interpreted): More than 6 hours on CPU
- Jags (interpreted): More than 128 GB, failed

# Compiler Engineering: a Principled Trade

*Some principles that can inform our work:*

- Compiler architecture
  - Pipeline of well-defined transformations and analysis
    - Closure conversion to implement first-class functions
    - Register allocation to use CPU registers
    - Abstract interpretation to cast static analysis
  - Between intermediate representations
    - Abstract syntax trees
    - Static Single Assignment
  - With well-defined semantics
- Compiler correctness
  - The compiler must preserve the semantics of the program
  - CompCert: a C compiler verified in Coq



# Compiler Engineering: a Principled Trade

*Some principles that can inform our work:*

- Compiler architecture
  - Pipeline of well-defined transformations and analysis
    - Closure conversion to implement first-class functions
    - Register allocation to use CPU registers
    - Abstract interpretation to cast static analysis
  - Between intermediate representations
    - Abstract syntax trees
    - Static Single Assignment
  - With well-defined semantics
- Compiler correctness
  - The compiler must preserve the semantics of the program
  - CompCert: a C compiler verified in Coq

# Compiler Engineering: a Principled Trade

*Some principles that can inform our work:*

- Compiler architecture
  - Pipeline of well-defined transformations and analysis
    - Closure conversion to implement first-class functions
    - Register allocation to use CPU registers
    - Abstract interpretation to cast static analysis
  - Between intermediate representations
    - Abstract syntax trees
    - Static Single Assignment
  - With well-defined semantics
- Compiler correctness
  - The compiler must preserve the semantics of the program
  - CompCert: a C compiler verified in Coq

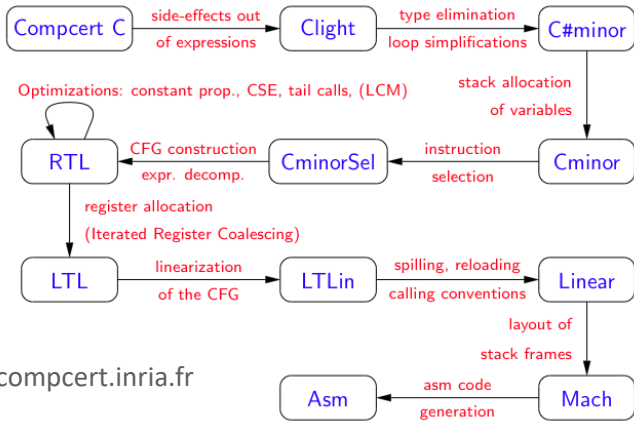
# Compiler Engineering: a Principled Trade

*Some principles that can inform our work:*

- Compiler architecture
  - Pipeline of well-defined transformations and analysis
    - Closure conversion to implement first-class functions
    - Register allocation to use CPU registers
    - Abstract interpretation to cast static analysis
  - Between intermediate representations
    - Abstract syntax trees
    - Static Single Assignment
  - With well-defined semantics
- Compiler correctness
  - The compiler must preserve the semantics of the program
  - CompCert: a C compiler verified in Coq

# Architecture of the Compcert Compiler

*Transformation is not one blob from C to Assembly*



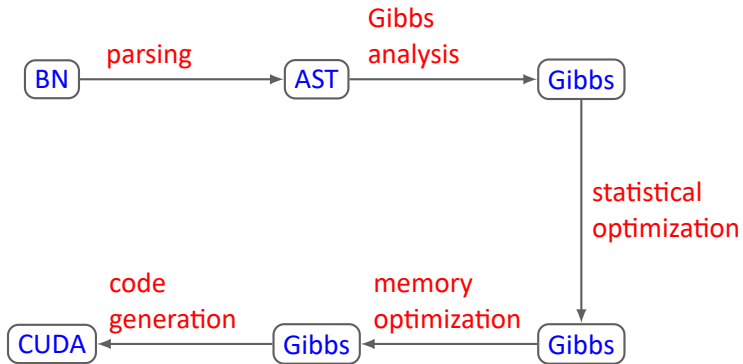
Source: [compcert.inria.fr](http://compcert.inria.fr)

# Compilation for Probabilistic Programs

These design principles can guide research in compilation:

- What are the key transformations?
  - *e.g.* How to derive a Gibbs sampler?
- How should we represent intermediate results?
  - Symbolic densities, “Gibbs equations”
- What are the high-level optimization specific to this field?
  - Approximation, estimation
- What is correctness for probabilistic program compilation?
  - *e.g.* Does compiler output converge to the query?

# Augur Compilation Pipeline

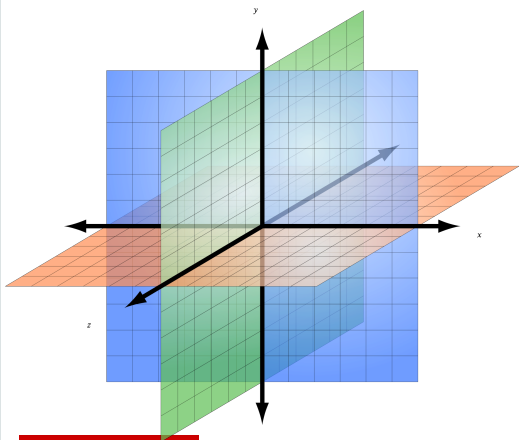


# Augur Compilation Pipeline

BN

# Language Classification

*We need to specify what language we're compiling*



3 dimensions of expressiveness

- ① Model (CRFs, BNs)
- ② Query (MAP, sampling)
- ③ Inference (MCMC, EM)



# Modeling Language

- Bayesian Networks
- Predefined family of distributions with known pdf/pmf
- Only bounded recursion, no conditionals, no functions

```
val phi = Dirichlet(V,beta).sample(K)
val theta = Dirichlet(K,alpha).sample(M)
val w = for(i <- 1 to M)
    for(j <- 1 to N(i))
        val z: Int = Categorical(K,theta(i)).sample()
        Categorical(V,phi(z)).sample()
```

# Query Language

- Fixed:

Maximum a Posteriori ( $\arg \max_{\theta, \phi} \mathbb{P}(\theta, \phi \mid \mathcal{W})$ )  
sampling

- Assign constants to subset of variables

```
observe(w, data)
```

# Inference

- Fixed:
  - Gibbs sampling
  - Metropolis-Hasting
  - Hamiltonian MC

We want to sample from  $p(X_1, X_2 | X_3 = x_3)$

- 1 Choose  $x_1^0, x_2^0$  at random
- 2 Repeat
  - $x_1^{t+1} \sim p(X_1 | X_2 = x_2^t, X_3 = x_3)$
  - $x_2^{t+1} \sim p(X_2 | X_1 = x_1^{t+1}, X_3 = x_3)$

These are called “Gibbs equations”

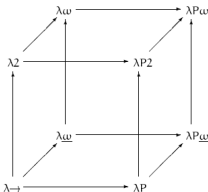
# We Want a Refund!

That's not *really* probabilistic programming!

Maybe, maybe not, but...

Can write *many, many* useful models

Incredibly useful to learn about probabilistic program compilation



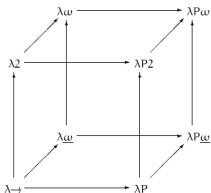
# We Want a Refund!

That's not *really* probabilistic programming!

Maybe, maybe not, but...

Can write *many, many* useful models

Incredibly useful to learn about probabilistic program compilation



# Augur Compilation Pipeline



# Representation: Abstract Syntax

How should we initially represent the program?

- We choose to symbolically represent the density

$$p ::= \text{cat} \mid \text{dir} \mid \text{bern}$$

$$P ::= \text{App}(p, \vec{X}) \mid \text{Cond}(p, \vec{X}, \vec{X}) \mid \text{Prod}(P, P) \mid \text{Mul}(i, N, P)$$

For our program, we have:

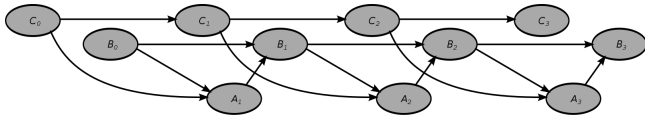
$$\text{Prod}(\text{Mul}(k = 1, K, \text{App}(\text{dir}, \phi_k)), \text{Mul}(m = 1, M, \text{Prod}(\dots, \dots)))$$

# Compiler vs Interpreter

Compiler:

$Prod(Mul(k = 1, K, App(dir, \phi_k)), Mul(m = 1, M, Prod(..., ...)))$

Interpreter:





# Semantics

- $\pi$ : abstract syntax of a probabilistic program
- $q$ : query (assignment of variables to values)

The semantics of  $\pi_q$  is a function  $\llbracket \pi_q \rrbracket$   
from unobserved variables to probabilities

In our case, we can define  $\llbracket \pi \rrbracket$  by induction on the AST

$$\begin{aligned}\llbracket \text{bern}(x, \mu) \rrbracket &= \mu^x (1 - \mu)^{(1-x)} \\ \llbracket \text{Prod}(t_1, t_2) \rrbracket &= \llbracket t_1 \rrbracket \times \llbracket t_2 \rrbracket\end{aligned}$$

# Semantics

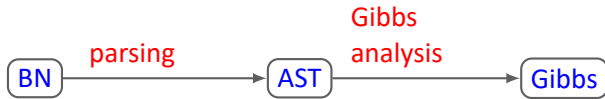
For LDA,  $\llbracket \pi \rrbracket =$

$$\lambda \phi, \theta, \mathbf{z}, \mathbf{w}. \left[ \prod_{k=1}^K \text{dir}(\phi_k) \right] \prod_{m=1}^M \text{dir}(\theta_m) \prod_{n=1}^{N_m} \text{cat}(w_{mn} | \phi_{z_{mn}}) \text{cat}(z_{mn} | \theta_m)$$

and

$$\llbracket \pi_{\mathbf{w}} \rrbracket = \lambda \phi. \lambda \theta. \lambda \mathbf{z} \frac{\llbracket \pi \rrbracket(\phi, \theta, \mathbf{z}, \mathbf{w})}{\sum_{\mathbf{w}} \llbracket \pi \rrbracket(\phi, \theta, \mathbf{z}, \mathbf{w})}$$

# Augur Compilation Pipeline



# Gibbs Sampling

We want to sample from  $p(X_1, X_2 | X_3 = x_3)$

- 1 Choose  $x_1^0, x_2^0$  at random
- 2 Repeat
  - $x_1^{t+1} \sim p(X_1 | X_2 = x_2^t, X_3 = x_3)$
  - $x_2^{t+1} \sim p(X_2 | X_1 = x_1^{t+1}, X_3 = x_3)$

These are called “Gibbs equations”

# Deriving Gibbs Equations by Rewriting

$$p(\phi_k | \mathbf{w}, \mathbf{z}, \boldsymbol{\theta}, \phi^{-k}) = ?$$

$$p(\theta_m | \mathbf{w}, \mathbf{z}, \boldsymbol{\theta}^{-m}, \phi) = ?$$

$$p(z_{mn} = k | w_{mn}, \theta_m, \phi_k) = ?$$

We derive the equations using a *term rewriting system*

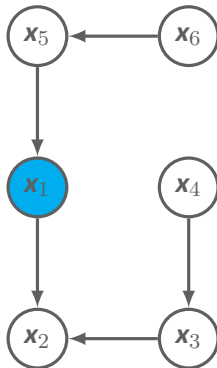
# Deriving Gibbs Equations by Rewriting

*The “Markov blanket” rule*

*Goal: get rid of conditionally independent variables*

Example rule:

$$P(x | y) \Rightarrow \frac{P(x,y)}{\int P(x,y) dx}$$



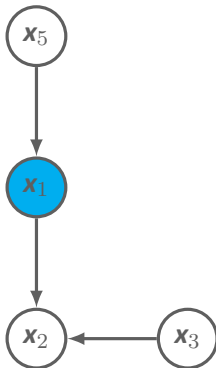
# Deriving Gibbs Equations by Rewriting

*The “Markov blanket” rule*

*Goal: get rid of conditionally independent variables*

Example rule:

$$P(x | y) \Rightarrow \frac{P(x,y)}{\int P(x,y) dx}$$

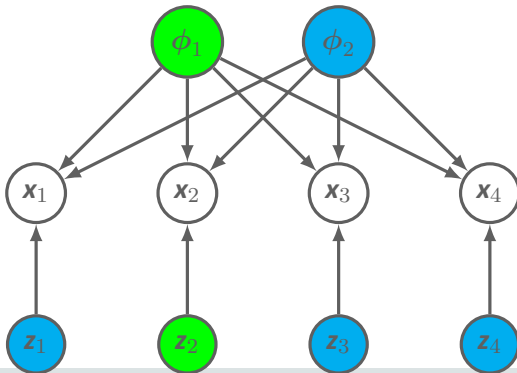


# Deriving Gibbs Equations by Rewriting

*The “mixture” rule*

*Goal: handling mixed and mixed-membership models*

$$\prod_i^N P(x_i) \Rightarrow \prod_i^N \{P(x_i)\}_{q(i)=T} \prod_i^N \{P(x_i)\}_{q(i)=F}$$



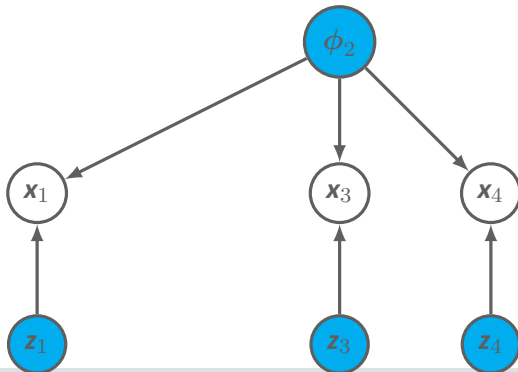


# Deriving Gibbs Equations by Rewriting

The “mixture” rule

Goal: handling mixed and mixed-membership models

$$\prod_i^N P(x_i) \Rightarrow \prod_i^N \{P(x_i)\}_{q(i)=T} \prod_i^N \{P(x_i)\}_{q(i)=F}$$



# Deriving Gibbs Equations by Rewriting

*The “conjugacy” rule*

*Goal: closed form solutions for some continuous variables*

Model: sequence  $x_1 \dots x_n$  of coin flips  $\prod_i \mu^{x_i} (1 - \mu)^{(1-x_i)}$

Query:  $p(\mu | x_1 \dots x_n)$ ?

$$p(\mu | x_1 \dots x_n) \Rightarrow p(\mu) \prod_i \mu^{x_i} (1 - \mu)^{(1-x_i)}$$
$$\Rightarrow ?$$

If  $p$  is Beta distribution then:

$$\Rightarrow \text{Beta}(\mu; \alpha + \sum_i \mathbb{1}[x_i = 1], \beta + \sum_i \mathbb{1}[x_i = 0])$$

# Deriving Gibbs Equations by Rewriting

*The “conjugacy” rule*

*Goal: closed form solutions for some continuous variables*

Model: sequence  $x_1 \dots x_n$  of coin flips  $\prod_i \mu^{x_i} (1 - \mu)^{(1-x_i)}$

Query:  $p(\mu | x_1 \dots x_n)$ ?

$$p(\mu | x_1 \dots x_n) \Rightarrow p(\mu) \prod_i \mu^{x_i} (1 - \mu)^{(1-x_i)}$$
$$\Rightarrow ?$$

If  $p$  is Beta distribution then:

$$\Rightarrow \text{Beta}(\mu; \alpha + \sum_i \mathbb{1}[x_i = 1], \beta + \sum_i \mathbb{1}[x_i = 0])$$

# Deriving Gibbs Equations by Rewriting

- Conjugate relation: sample from known distribution
- Non conjugate relation: sample using, *e.g.* rejection sampling
- Discovering these relations is critical
- Pointing out when they are available is useful
  - Something an IDE could do

# The Gibbs Equations

$$p(\phi_k | \mathbf{w}, \mathbf{z}, \boldsymbol{\theta}, \phi^{-k}) = \text{dir}(\phi_k | \langle \text{wpt}[k][v] + \beta : v \in [1..V] \rangle)$$

$$p(\theta_m | \mathbf{w}, \mathbf{z}, \boldsymbol{\theta}^{-m}, \phi) = \text{dir}(\theta_m | \langle \text{tpd}[m][k] + \alpha : k \in [1..K] \rangle)$$

$$p(z_{mn} = k | w_{mn}, \theta_m, \phi_k) = \text{cat}(w_{mn} | \phi_k) \text{cat}(k | \theta_m)$$

$\text{tpd}[m][k]$ : count occurrences of topic  $k$  in document  $m$

$\text{wpt}[k][n]$ : count occurrences of word  $n$  assigned to topic  $k$

Independence  $\Rightarrow$  parallelism

# The Final Code

In parallel:  $\phi_k \sim \text{Dir}(\langle \text{wpt}[k][v] + \beta : v \in [1..V] \rangle)$

In parallel:  $\theta_m \sim \text{Dir}(\langle \text{tpd}[m][k] + \alpha : k \in [1..K] \rangle)$

In parallel:  $z_{mn} \sim \phi'_{w_{mn}} \odot \theta_m$

$\text{tpd}[m][k]$ : count occurrences of topic  $k$  in document  $m$

$\text{wpt}[k][n]$ : count occurrences of word  $n$  assigned to topic  $k$

# Compiler vs Interpreter

## Compiler:

### Term rewriting

- Difficult to do, balancing ad hoc and impossible
- Does the rewrite process terminate?
- Can we characterize when it fails?
- + Having the equations is highly valuable
- + Can be parallelized very effectively for GPU

## Interpreter:

At runtime, for each variable, “read off” the blanket

- + Trivial to do
- Expensive, “local” view of the program

# Compiler Correctness

*What would it mean for this compiler to be correct?*

- Program  $\pi$  + query  $q$
- Initial state sampled from  $\mu$
- The stochastic matrix for  $\pi_q$  is  $P_{\pi_q}$
- $\delta$ : total variation distance

$$\forall \pi_q, \exists \tau, \forall \mu, \delta(\mu P_{\pi_q}^t, \tau) \rightarrow 0 \text{ as } t \rightarrow \infty$$

Not sufficient.



# Compiler Correctness

*What would it mean for this compiler to be correct?*

- Program  $\pi$  + query  $q$
- Initial state sampled from  $\mu$
- The stochastic matrix for  $\pi_q$  is  $P_{\pi_q}$
- $\delta$ : total variation distance

$$\forall \pi_q, \exists \tau, \forall \mu, \delta(\mu P_{\pi_q}^t, \tau) \rightarrow 0 \text{ as } t \rightarrow \infty$$

Not sufficient.

# Compiler Correctness

- Program  $\pi$  + query  $q$
- Initial state sampled from  $\mu$
- The stochastic matrix for  $\pi_q$  is  $P_{\pi_q}$
- $\delta$ : total variation distance

$$\forall \pi_q, \delta(\mu P_{\pi_q}^t, \llbracket \pi_q \rrbracket) \rightarrow 0 \text{ as } t \rightarrow \infty$$

Sufficient. Unfortunately, not necessary... but this is another talk

# Compiler Correctness

- Program  $\pi$  + query  $q$
- Initial state sampled from  $\mu$
- The stochastic matrix for  $\pi_q$  is  $P_{\pi_q}$
- $\delta$ : total variation distance

$$\forall \pi_q, \delta(\mu P_{\pi_q}^t, \llbracket \pi_q \rrbracket) \rightarrow 0 \text{ as } t \rightarrow \infty$$

Sufficient. Unfortunately, not necessary... but this is another talk

# Augur: Conclusion

- Some successes
  - We can generate very efficient distributed/GPU code for non-trivial models/inference
  - It's important to carefully represent the intermediate program
- Some remaining challenges
  - The term rewriting system is very unprincipled
  - Compiler code-base explosion to support more inference methods

Check out this work: [Hakaru](#), [Shuffle](#), [Bayonet](#)

# Augur: Conclusion

- Some successes
  - We can generate very efficient distributed/GPU code for non-trivial models/inference
  - It's important to carefully represent the intermediate program
- Some remaining challenges
  - The term rewriting system is very unprincipled
  - Compiler code-base explosion to support more inference methods

Check out this work: Hakaru, Shuffle, Bayonet

# Conclusion

For machine learning people:

Think in terms of:

- Narrowly-scoped transformations
- Intermediate representations
- Transformation correctness

For compiler people:

- Separation of concerns more difficult for probabilistic programming
- Some transformations (like Gibbs) deserve new ideas
- Verification milestone: prove correctness formally

# More Info

## Collaborators:

Oracle Labs ML: **Adam Pockock**, Steve Green

Oracle Labs PL: Guy Steele

Harvard: **Daniel Huang**, Greg Morrisett

CMU: Joseph Tassarotti

## Papers:

NIPS'14 Augur: Data-Parallel Probabilistic Modeling

PLDI'17 Compiling Markov chain Monte Carlo algorithms for probabilistic modeling

Code: Augur v2, by Daniel Huang

- <https://github.com/danehuang/augurv2>

# Integrated Cloud

## Applications & Platform Services



ORACLE®