

AI-Systems Machine Learning Frameworks

Joseph E. Gonzalez

Co-director of the RISE Lab

jegonzal@cs.berkeley.edu

Objectives For Today

- Historical Evolution of Machine Learning Frameworks
- Declarative (Lazy) vs Imperative (Eager) DSLs
- Automatic Differentiation
- This weeks reading

Historical Context

Early ML / Stats Languages

- **S Data Programming** Languages
 - Developed in **1976** as **Bell Labs** by John Chambers
 - **Replaced Fortran** by providing higher level APIs, graphics
 - Developed **formula syntax** for describing models
 - Eventually replaced by R ...
- **R open-source** implementation of S (S-Plus)
 - Developed in **1990's** at University of Auckland
 - Ross Ihaka, Robert Gentleman
 - Like S/S-Plus → Linear algebra abstractions
 - **Rich set of libraries** for statistical analysis
 - Still widely used

- **Matlab** (Matrix Laboratory) – Numerical Computing Sys.
 - Developed in **1970s** at the University of New Mexico by Cleve Moler
 - Designed to **simplify access** to **LINPACK** and **EISPACK**
 - Reasonable **integration with C/Fortran**
 - Rich **graphical interface** with support for graphical programming
 - Simulink
 - Expensive → Octave **limited** open-source version
 - Popular in applied math, engineering, and controls community
 - **Extremely popular in the machine learning community**
 - We would joke that ML people only knew how to program Matlab
- and then it all changed ...

Rise of the **Python** Eco-System

- Development of **%pylab**
 - iPython (2001) + SciPy (2001) + Matplotlib (2003) + NumPy (2006)
 - Functions /APIs were like Matlab so easy to transition
 - **Freeeeeeee!**
- **Scikit-learn** – basic ML algorithms and models (2007)
 - Started as Google summer of code project → **developed by INRIA**
 - Wide range of standard machine learning techniques
- ~2012 large fraction of ML community Matlab → Python
 - Why?
- Development remained focused on **algorithms libraries**

Machine Learning Libraries

- **LIBLINEAR/LIBSVM** (2008) – **fast algorithms** for fitting linear models and kernelized SVMs
 - Developed at National Taiwan University for (still used in Sklearn)
- **Vowpal Wabbit** (2010?) – **out-of-core** learning for generalized linear models and others
 - Developed by John Langford while at Yahoo!
 - Popular for high-dimensional features
- **Weka** (Java version 1997) – Collection of ML algorithms for Java
 - Developed at the University of Waikato in New Zealand
 - Provided tools for visualizing and analyzing data
- **Xgboost** (2014) – **distributed** boosted decision trees
 - Developed by Tianqi Chen at University of Washington
- Many more ...

Distributed Machine Learning Frameworks

- **Mahout** (2009) – **ML algorithms** on Hadoop
 - Early distributed ML library with “**recommender algorithms**”
 - Unable to leverage memory caching
- **GraphLab** (2010) – Framework for **graph structured algorithms**
 - Contained library of algs. (e.g., Gibbs Sampling, LoopyBP, ...)
 - Developed new abstractions for distributed graph algs.
- **Spark mllib / SparkML** (2014) – ML algorithms for Spark
 - Leverages memory caching
 - Benefits from work on GraphLab/Sklearn/SystemML

Languages vs Algorithm Libraries



- **Languages** provided support for mathematical operations
 - User still implemented new models and algorithms using fundamental **linear algebra primitives**
- **Libraries of Algorithms** provided individual learning techniques
 - Often specialized to model/technique (fast and easy-to-use)
- Need something in the middle!

Embedded Domain Specific Languages

- Domain specific languages (DSLs) provide **specialized functionality** for a given task
 - Limited functionality → **simplicity** and **optimization**
 - **Example:** SQL → Specialized for data manipulation
- Embedded DSLs are **libraries** or **language extensions** within a general-purpose language tailored to a specific task
 - Combine benefits of DSL and general languages
 - **Example:** linear algebra libraries
- Embedded DSLs have played a significant role in ML
 - Linear Algebra → Pipelines → Differentiable Programs

Machine Learning Pipelines

- Scikit Learn Pipelines (2011)
 - Describes **composition** of feature transformations and models
 - Enables **end-to-end training** and **standardized** prediction

```
steps = [('scaler', StandardScaler()), ('SVM', SVC())]  
pipeline = Pipeline(steps) # define the pipeline object.  
parameteres = {'SVM__C':[0.001,0.1,10,100,10e5], 'SVM__gamma':[0.1,0.01]}  
grid = GridSearchCV(pipeline, param_grid=parameteres, cv=5)  
grid.fit(X_train, y_train)
```

- Spark ML Pipelines (Similar to SkLearn)

```
tokenizer = Tokenizer(inputCol="text", outputCol="words")  
hashingTF = HashingTF(inputCol=tokenizer.getOutputCol(), outputCol="features")  
lr = LogisticRegression(maxIter=10, regParam=0.001)  
pipeline = Pipeline(stages=[tokenizer, hashingTF, lr])
```

SystemML (VLDB'16)

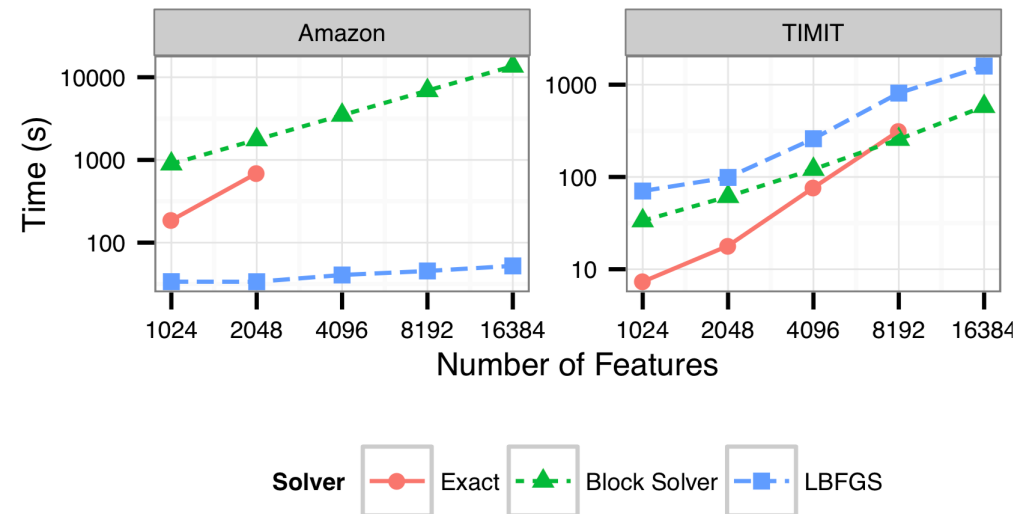
- Developed at IBM
- **Domain specific language** for describing ML algorithms
 - Python/R like but **not embedded**
 - Optimizer and runtime to execute on **Apache Spark**
- Explored range of optimizations
 - Data repartitioning
 - Caching
 - Distributed matrix representations

```
1: X = read($inFile);
2: r = $rank; lambda = $lambda;
3: U = rand(rows=nrow(X), cols=r, min=-1.0, max=1.0);
4: V = rand(rows=r, cols=ncol(X), min=-1.0, max=1.0);
5: W = (X != 0);
6: mi = $maxiter; mii = r; i = 0; is_U = TRUE;
7: while(i < mi) {
8:   i = i + 1; ii = 1;
9:   if (is_U)
10:     G = (W * (U %*% V - X)) %*% t(V) + lambda * U;
11:   else
12:     G = t(U) %*% (W * (U %*% V - X)) + lambda * V;
13:   norm_G2 = sum(G ^ 2); norm_R2 = norm_G2;
14:   R = -G; S = R;
15:   while(norm_R2 > 10E-9 * norm_G2 & ii <= mii) {
16:     if (is_U) {
17:       HS = (W * (S %*% V)) %*% t(V) + lambda * S;
18:       alpha = norm_R2 / sum (S * HS);
19:       U = U + alpha * S;
20:     } else {
21:       HS = t(U) %*% (W * (U %*% S)) + lambda * S;
22:       alpha = norm_R2 / sum (S * HS);
23:       V = V + alpha * S;
24:     }
25:     R = R - alpha * HS;
26:     old_norm_R2 = norm_R2; norm_R2 = sum(R ^ 2);
27:     S = R + (norm_R2 / old_norm_R2) * S;
28:     ii = ii + 1;
29:   }
30:   is_U = ! is_U;
31: }
32: write(U, $outUFile, format = "text");
33: write(V, $outVFile, format = "text");
```

Keystone ML (ICDE'17)

- Developed in AMPLab@Berkeley
- Pipelines of **ML algorithms** and optimization on top of Spark
 - **Embedded Scala DSL**
 - Outperformed SystemML
- Cost based optimize to select best version of learning algorithm based on inputs
 - Example: QR vs L-BFGS

```
val textClassifier = Trim andThen  
  LowerCase andThen  
  Tokenizer andThen  
  NGramsFeaturizer(1 to 2) andThen  
  TermFrequency(x => 1) andThen  
  (CommonSparseFeatures(1e5), data) andThen  
  (LinearSolver(), data, labels)  
val predictions = textClassifier(testData)
```



Languages vs Algorithm Libraries



- Increased focus on deep learning → empirical risk minimization for complex **differentiable models**
- Research shifts from algorithm design to **model design**
- **Deep Learning Frameworks:** Theano (2008), Caffe (2014), MXNet (2015), TensorFlow (2015), PyTorch (2016)
 - Combine **automatic differentiation** with **hardware acceleration**

Review of Automatic Differentiation

Automatic Differentiation

- Method of computing **numeric derivatives** of a **program** by **tracking** the **forward execution** of that **program**
- Other methods for computing derivatives
 - **Manual implementation:** the standard method in deep learning prior to these frameworks
 - laborious and **error prone!**
 - **Numerical differentiation:** using finite differences
 - Easy, costly and sensitive to numerical precision
 - **Symbolic differentiation:** using computer algebraic systems
 - Expressions can grow exponentially

Illustration from “Automatic Differentiation in Machine Learning: a Survey”

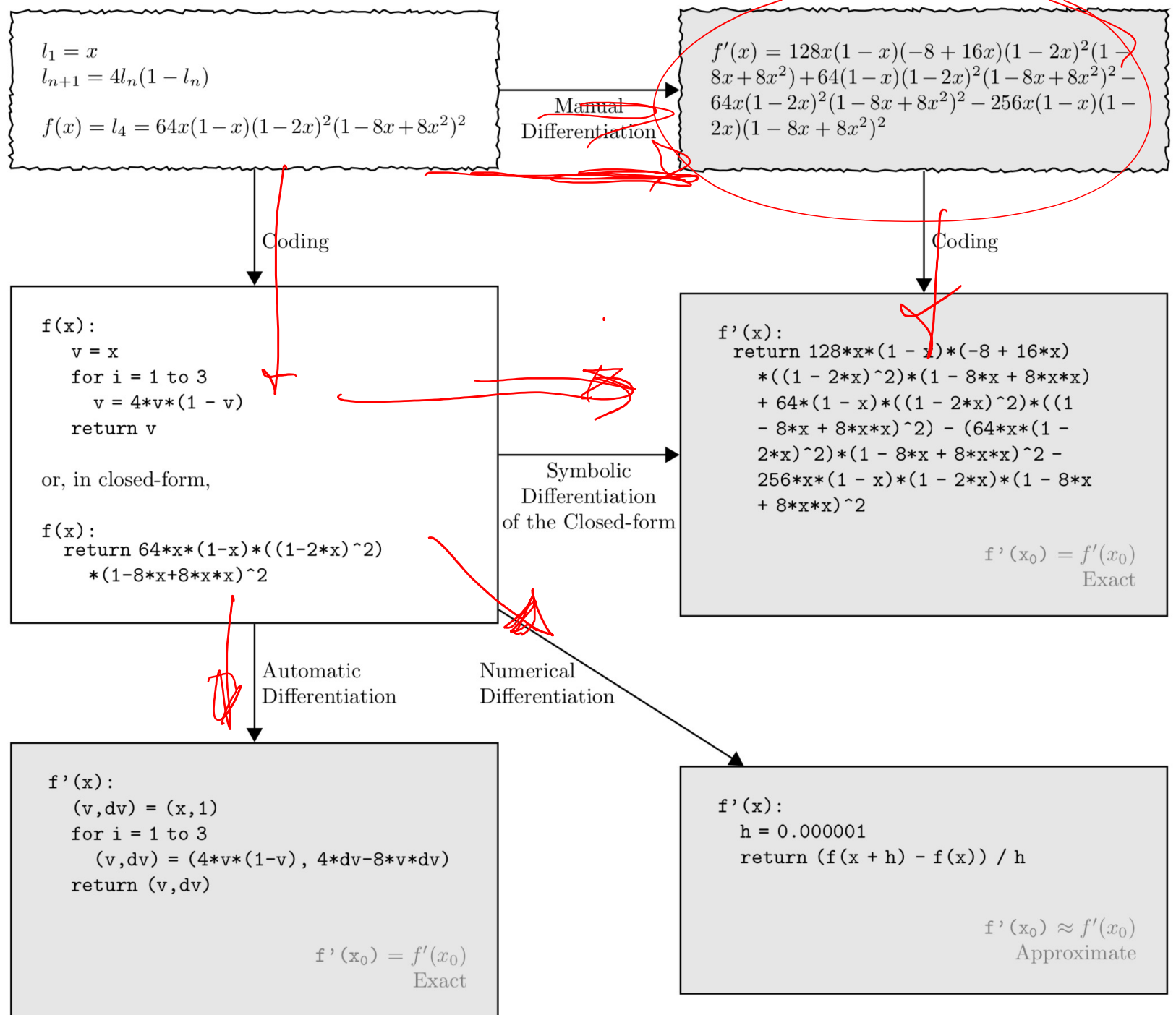
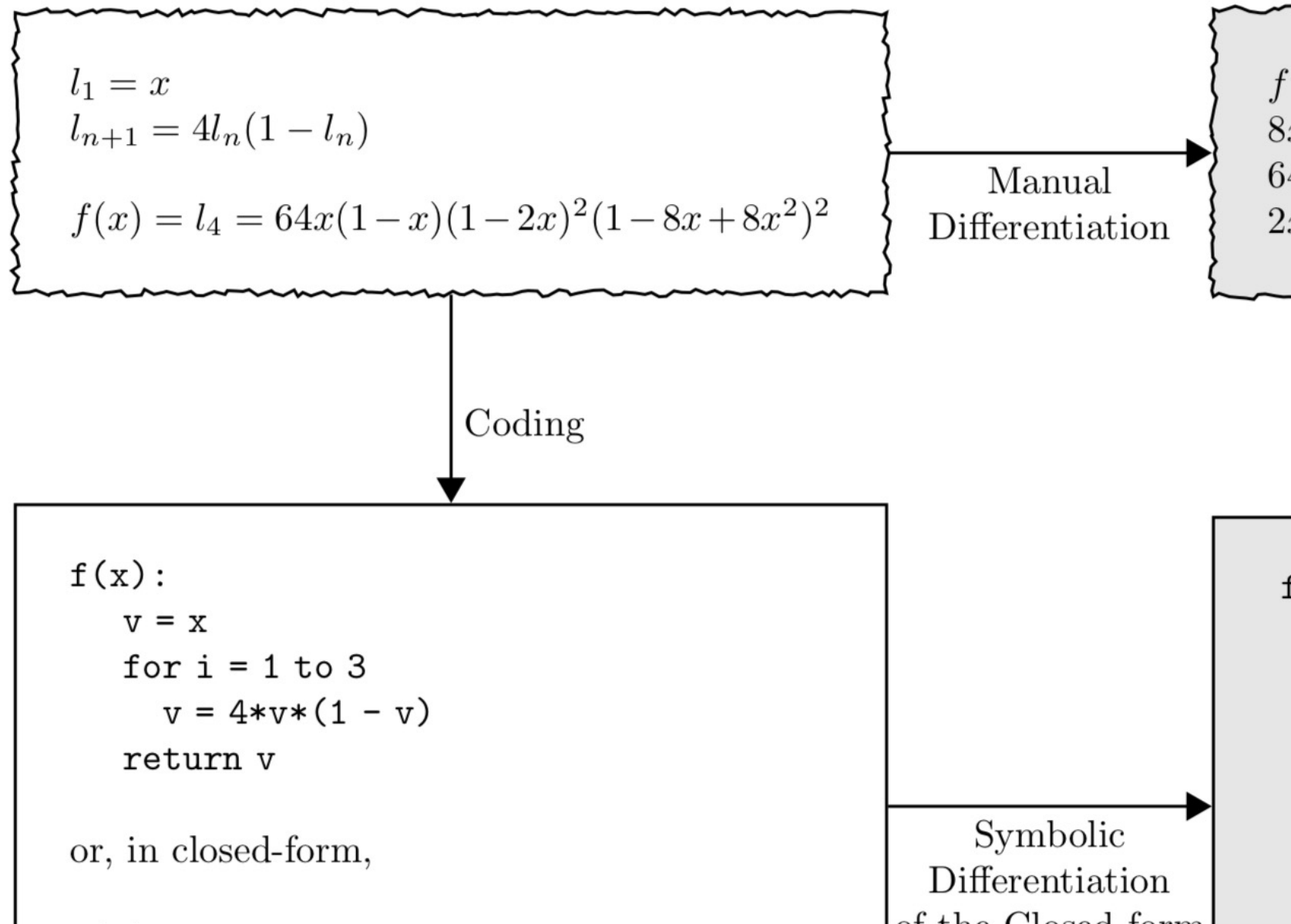


Illustration from
“Automatic
Differentiation in
Machine Learning:
a Survey”



$$f(x) = x(1-x)(1-2x)^2(1-8x+8x^2)^2$$

Manual
Differentiation

$$f'(x) = 128x(1-x)(-8+16x)(1-2x)^2(1-8x+8x^2) + 64(1-x)(1-2x)^2(1-8x+8x^2)^2 - 64x(1-2x)^2(1-8x+8x^2)^2 - 256x(1-x)(1-2x)(1-8x+8x^2)^2$$

Coding

to 3
(1 - v)
term,

Coding

```
f'(x):
return 128*x*(1-x)*(-8+16*x)
      *((1-2*x)^2)*(1-8*x+8*x*x)
      + 64*(1-x)*((1-2*x)^2)*((1-8*x+8*x*x)^2) - (64*x*(1-2*x)^2)*(1-8*x+8*x*x)^2 - 256*x*(1-x)*(1-2*x)*(1-8*x+8*x*x)^2
```

Symbolic
Differentiation

$$l_1 = x$$

$$l_{n+1} = 4l_n(1 - l_n)$$

$$f(x) = l_4 = 64x(1-x)(1-2x)^2(1-8x+8x^2)^2$$

Manual
Differentiation

$$f'(x) = 128x(1-x)(-8+16x)(1-2x)^2(1-8x+8x^2) + 64(1-x)(1-2x)^2(1-8x+8x^2)^2 - 64x(1-2x)^2(1-8x+8x^2)^2 - 256x(1-x)(1-2x)(1-8x+8x^2)^2$$

Coding

```
f(x):
  v = x
  for i = 1 to 3
    v = 4*v*(1 - v)
  return v
```

or, in closed-form,

```
f(x):
  return 64*x*(1-x)*((1-2*x)^2)
  *(1-8*x+8*x*x)^2
```

Symbolic
Differentiation
of the Closed-form

```
f'(x):
  return 128*x*(1-x)*(-8+16*x)
  *((1-2*x)^2)*(1-8*x+8*x*x)
  + 64*(1-x)*((1-2*x)^2)*((1-8*x+8*x*x)^2)
  - (64*x*(1-2*x)^2)*(1-8*x+8*x*x)^2
  - 256*x*(1-x)*(1-2*x)*(1-8*x+8*x*x)^2
```

$$f'(x_0) = f'(x_0)$$

Exact

How I used to do this as a
graduate student (2010).

How I would cheat using
Mathematica.

$$l_1 = x$$

$$l_{n+1} = 4l_n(1 - l_n)$$

$$f(x) = l_4 = 64x(1-x)(1-2x)^2(1-8x+8x^2)^2$$

Manual
Differentiation

$$f'(x) = 128x(1-x)(-8+16x)(1-2x)^2(1-8x+8x^2) + 64(1-x)(1-2x)^2(1-8x+8x^2)^2 - 64x(1-2x)^2(1-8x+8x^2)^2 - 256x(1-x)(1-2x)(1-8x+8x^2)^2$$

Coding

```
f(x):
    v = x
    for i = 1 to 3
        v = 4*v*(1 - v)
    return v
```

or, in closed-form,

```
f(x):
    return 64*x*(1-x)*((1-2*x)^2)
        *(1-8*x+8*x*x)^2
```

Automatic
Differentiation

```
f'(x):
    (v,dv) = (x,1)
    for i = 1 to 3
        (v,dv) = (4*v*(1-v), 4*dv-8*v*dv)
    return (v,dv)
```

Coding

```
f'(x):
    return 128*x*(1-x)*(-8+16*x)
        *((1-2*x)^2)*(1-8*x+8*x*x)
        + 64*(1-x)*(1-2*x)^2*(1-8*x+8*x*x)^2
        - 64*x*(1-2*x)^2*(1-8*x+8*x*x)^2
        - 256*x*(1-x)*(1-2*x)*(1-8*x+8*x*x)^2
```

Symbolic
Differentiation
of the Closed-Form

Numerical
Differentiation

```
f'(x):
    h = 0.000001
    return (f(x+h) - f(x)) / h
```

$$f'(x_0) \approx \frac{f(x_0+h) - f(x_0)}{h}$$

Automatic differentiation
operates on a program to
generate a program that
computes the derivative
efficiently and accurately.

Key Ideas in Automatic Differentiation

- Leverage **Chain Rule** to reason about function composition

$$\frac{\partial}{\partial x} f(g(x)) = \dot{f}(g(x)) \frac{\partial}{\partial x} g(x)$$

- Two modes of automatic differentiation
 - **Forward differentiation**: computes derivative during execution
 - efficient for single derivative with multiple outputs
 - **Backward differentiation (back-propagation)**: computes derivative (gradient) by reverse evaluation of the computation graph
 - Efficient for multiple derivative (gradient) calculation + **Requires caching**

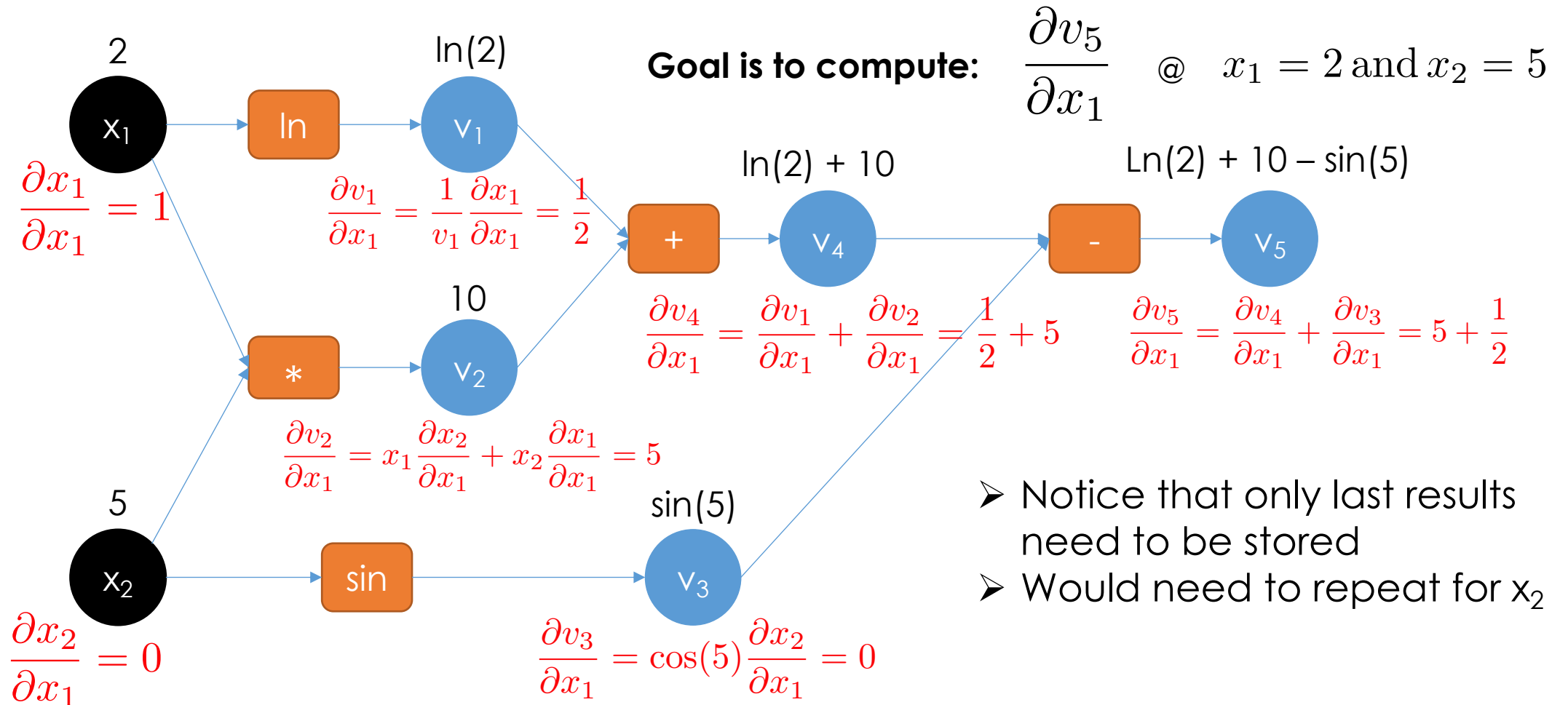
Forward Differentiation (Example)

$$f(x_1, x_2) = \ln(x_1) + x_1 x_2 - \sin(x_2)$$

Goal is to compute: $\frac{\partial v_5}{\partial x_1}$ @ $x_1 = 2$ and $x_2 = 5$

Forward Differentiation (Example)

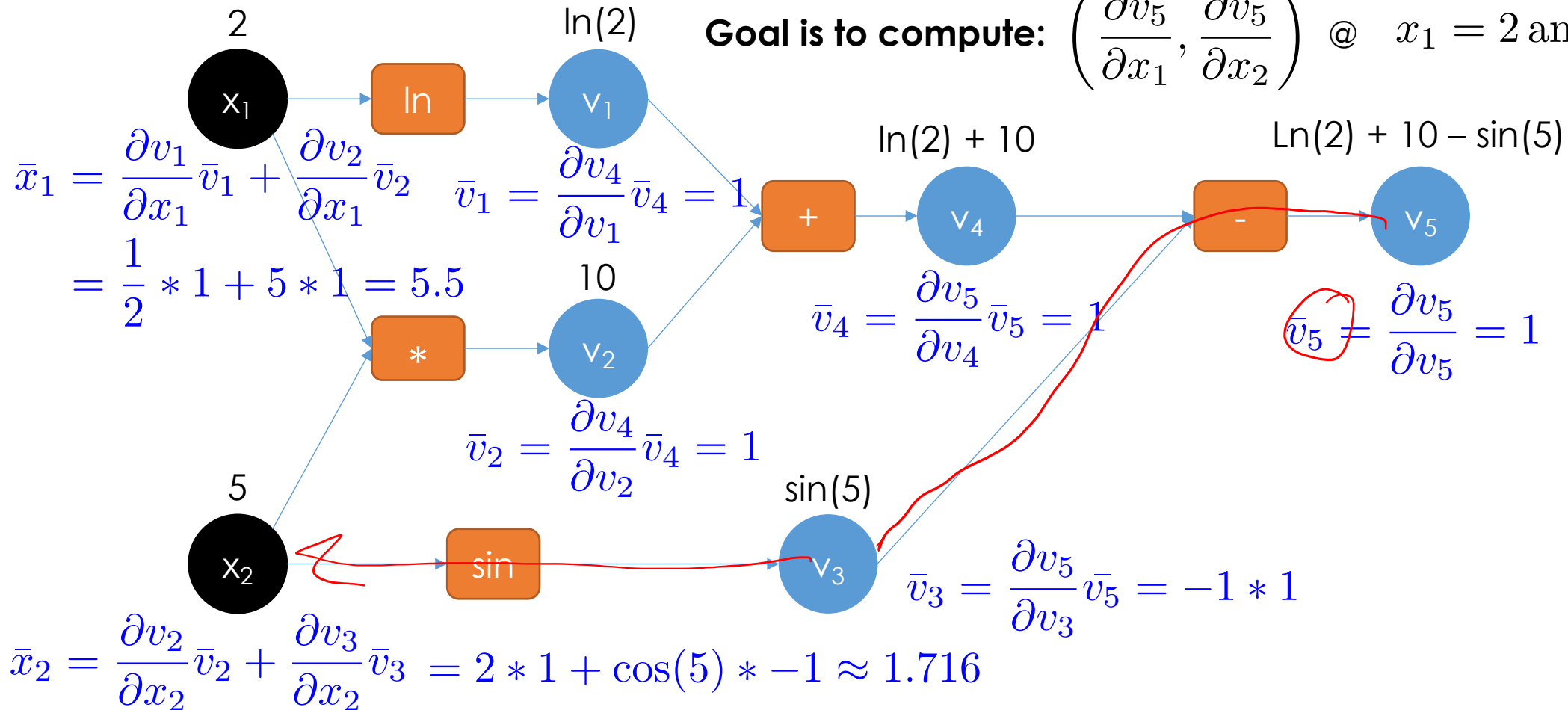
$$f(x_1, x_2) = \ln(x_1) + x_1 x_2 - \sin(x_2)$$



Backward (Reverse) Differentiation

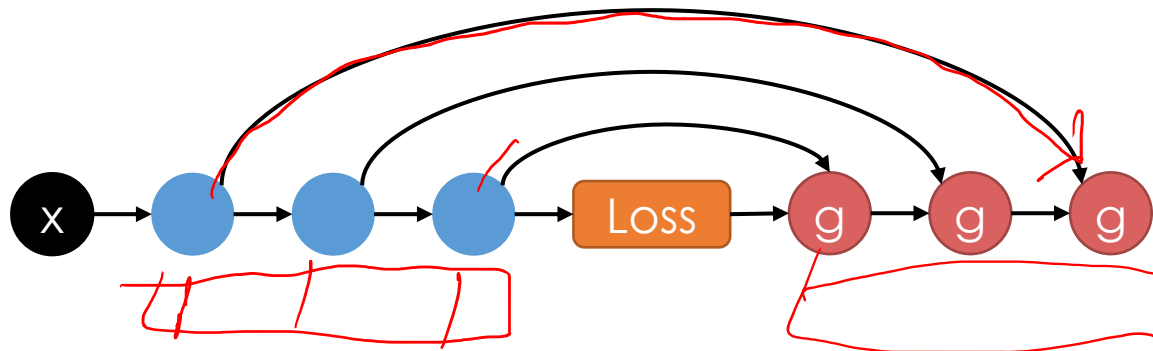
$$f(x_1, x_2) = \ln(x_1) + x_1 x_2 - \sin(x_2)$$

Goal is to compute: $\left(\frac{\partial v_5}{\partial x_1}, \frac{\partial v_5}{\partial x_2} \right) @ x_1 = 2 \text{ and } x_2 = 5$



Backward (Reverse) Differentiation

- Performs well when **computing large gradients** relative to number of function outputs
 - When might forward differentiation perform well? Why?
- Requires **caching** or **recomputing** intermediate activations from forward pass
 - Active research on what to recompute vs cache



Deep Learning Frameworks

Declarative vs Imperative Abstractions

- **Declarative** (*define-and-run*): Embedded DSL used to construct **static computation graph**
 - Examples: Theano (2010), Caffe (2014), TensorFlow (2015)
 - **Easier** to optimize, distribute, and export models
- **Imperative** (*define-by-run*): Embedded DSL used to directly compute output resulting in a **dynamic computation graph** defined by the program
 - Examples: Chainer (2015), autograd (2016), PyTorch (2017)
 - **Interpreted execution** of inference and gradient
 - **Easier** to program and debug
- **Hybrid Approaches**: Current research
 - TensorFlow Eager, MXNet

Theano – Original Deep Learning Framework

- First developed at the **University of Montreal** (2008)
 - from **Yoshua Bengio's** group
- **Abstraction:** Python embedded DSL (as a library) to construct symbolic expression graphs for complex mathematical expressions
- **System:** a **compiler** for **mathematical expressions** in Python
 - Optimizes mathematical expressions (e.g., $(A+b)(A+b)=(A+b)^2$)
 - CPU/**GPU acceleration**
 - Also ... **automatic differentiation**

```
import numpy
import theano.tensor as T
from theano import shared, function
```

Declaring Variables

```
x = T.matrix()
y = T.lvector()
w = shared(numpy.random.randn(100))
b = shared(numpy.zeros(()))
```

```
print "Initial model:"
print w.get_value(), b.get_value()
```

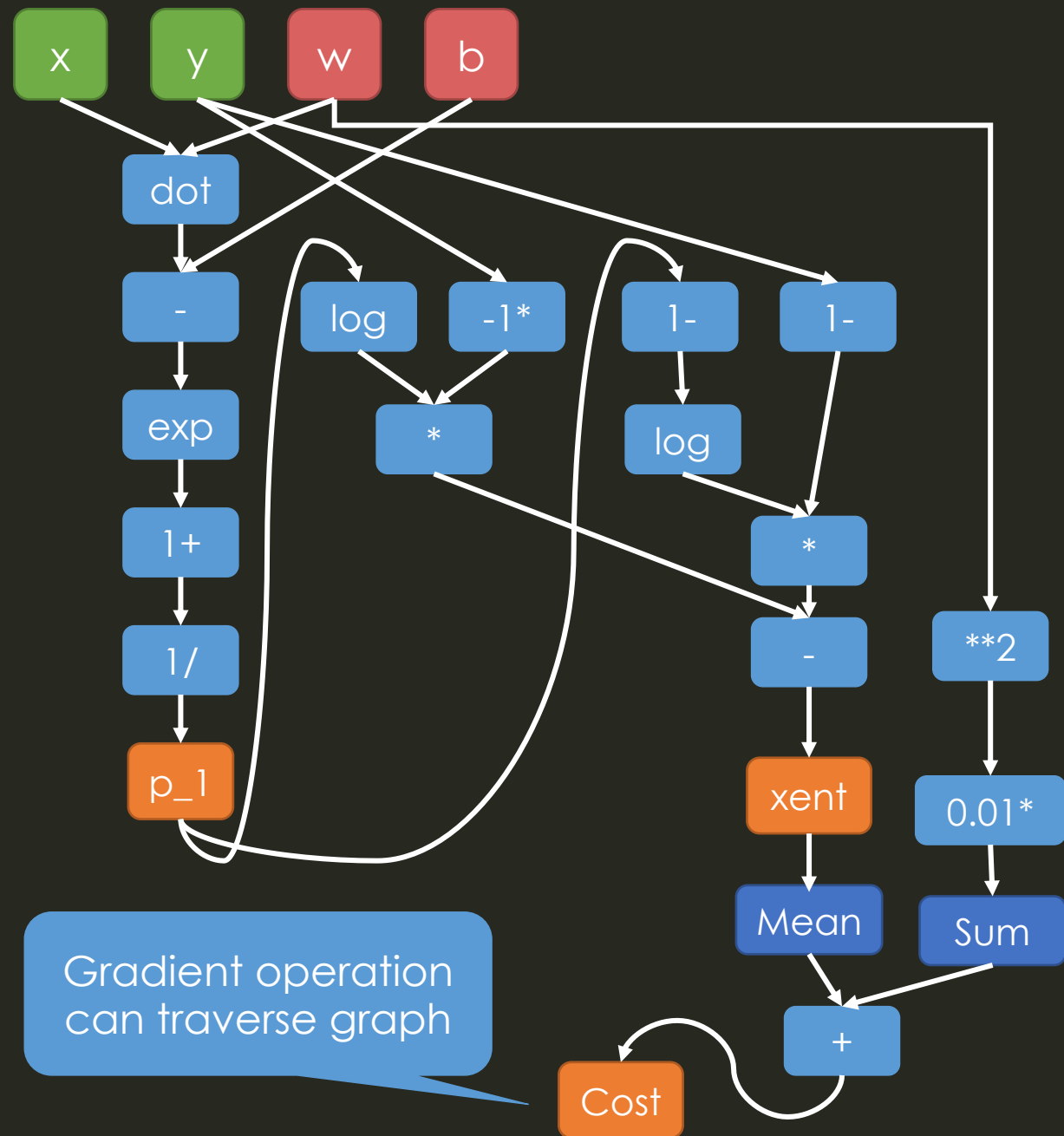
```
p_1 = 1 / (1 + T.exp(-T.dot(x, w)-b))
xent = -y*T.log(p_1) - (1-y)*T.log(1-p_1)
cost = xent.mean() + 0.01*(w**2).sum()
gw,gb = T.grad(cost, [w,b])
prediction = p_1 > 0.5
```

What is the value (type) of prediction?

Building Expression Graph

Note that this looks like a NumPy expression

This is more difficult to debug and reason about.



```
import numpy
import theano.tensor as T
from theano import shared, function
```

Declaring Variables

```
x = T.matrix()
y = T.lvector()
w = shared(numpy.random.randn(100))
b = shared(numpy.zeros(()))
```

```
print "Initial model:"
print w.get_value(), b.get_value()
```

```
p_1 = 1 / (1 + T.exp(-T.dot(x, w)-b))
xent = -y*T.log(p_1) - (1-y)*T.log(1-p_1)
cost = xent.mean() + 0.01*(w**2).sum()
gw,gb = T.grad(cost, [w,b])
prediction = p_1 > 0.5
```

What is the value (type) of prediction?

Building Expression Graph

Note that this looks like a NumPy expression

This is more difficult to debug and reason about.

```
predict = function(inputs=[x],
                   outputs=prediction)
train = function(
    inputs=[x,y],
    outputs=[prediction, xent],
    updates={w: w - 0.1*gw, b: b - 0.1*gb})
```

Updates shared variables after computation

```
N = 4
feats = 100
D = (numpy.random.randn(N, feats),
     numpy.random.randint(size=N, low=0, high=2))
training_steps = 10
for i in range(training_steps):
    pred, err = train(D[0], D[1])
print "Final model:",
print w.get_value(), b.get_value()
print "target values for D", D[1]
print "prediction on D", predict(D[0])
```

Instantiating Values

Function call compiles graphs into optimized native execution.

Theano Compilation of Functions



- **Rewriting** (simplify) mathematical expression
 - $\text{Exp}(\log(x)) = x$
- Duplicate **code elimination**
 - Important because gradient rewrites introduce redundancy
 - Recall gradient calculations extend graph via the chain rule

Theano Compilation of Functions



Addresses **numerical stability** of operations

➤ Example: for $x = 709$, $x = 710$ what is the value of

$$\log(1 + \exp(x)) =$$

- for $x = 709 \rightarrow 709$
- for $x = 710 \rightarrow \text{inf}$
- Rewritten as x for $x > 709$

Theano Compilation of Functions



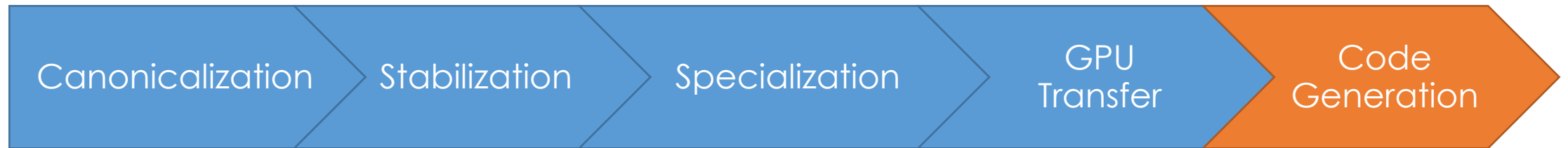
- Rewrite subgraphs to more efficient forms
 - `pow(x,2) → square(x)`
 - Tensor slicing → memory aliasing
 - Mapping to best version of GEMM routines

Theano Compilation of Functions



- GPU versions of ops are introduced (where possible)
- Copy routines are added to move data

Theano Compilation of Functions



- Generate and link C++ and CUDA implementations of operators
 - Picking from existing implementations
 - Specialization for different dtypes

What happened to Theano?

- Fairly advanced compared to TensorFlow (TF) in 2016
 - Symbolic gradient optimization and wide range of operators
 - Initially **faster than TensorFlow**
- What happened? (some speculation...)
 - Didn't have the backing of a large industrial group
 - TensorFlow was being pushed heavily by Google
 - Did not support multi-GPU/distributed computation and limited support for user defined parallelization
 - TensorFlow had more built-in deep learning operators
 - Theano lacked visualization tools (e.g., TensorBoard)
 - Complaints about error messages...?

PyTorch

- **Imperative DL library** which works like NumPy (on GPUs)

```
if torch.cuda.is_available():
    device = torch.device("cuda") # a CUDA device object
    y = torch.ones_like(x, device=device) # directly create a tensor on GPU
    x = x.to(device) # or just use strings ``.to("cuda")``
    z = x + y
    print(z) # tensor([2.0814], device='cuda:0')
    print(z.to("cpu", torch.double)) # tensor([2.0814], dtype=torch.float64) er!
```

- and supports automatic differentiation

```
x = torch.ones(2, 2, requires_grad=True) # tensor([[3., 3.],
y = x + 2 # [3., 3.], grad_fn=<AddBackward0>)
print(y)
z = y * y * 3 # tensor(27., grad_fn=<MeanBackward0>)
out = z.mean()
print(out) # tensor([4.5000, 4.5000],
out.backward() # [4.5000, 4.5000]])
print(x.grad)
```

This weeks readings

Reading for the Week

- [Automatic differentiation in ML: Where we are and where we should be going](#)
 - NeurIPS'18
 - Provides an overview of the state of automatic differentiation
- [TensorFlow: A System for Large-Scale Machine Learning](#)
 - OSDI'16
 - The primary TensorFlow paper discusses system and design goals
- [JANUS: Fast and Flexible Deep Learning via Symbolic Graph Execution of Imperative Programs](#)
 - NSDI'19
 - Recent work exploring a method to bridge Declarative and Imperative approaches in TensorFlow

Extra Suggested Reading

- [Automatic Differentiation in Machine Learning: a Survey](#)(JMLR'18)
 - Longer discussion on automatic differentiation in ML
- [Theano: A CPU and GPU Math Compiler in Python](#)(SciPy'10)
 - Great overview of AD and Theano system
- [TensorFlow Eager: A Multi-Stage, Python-Embedded DSL for Machine Learning](#) (arXiv'19)
 - Good follow-up to TF paper addressing limitations

Automatic differentiation in ML: Where we are and where we should be going?

Bart van Merriënboer, Olivier Breuleux, Arnaud Bergeron,
Pascal Lamblin

From Mila (home of Theano) and Google Brain (home of TF)

Automatic differentiation in ML: Where we are and where we should be going?

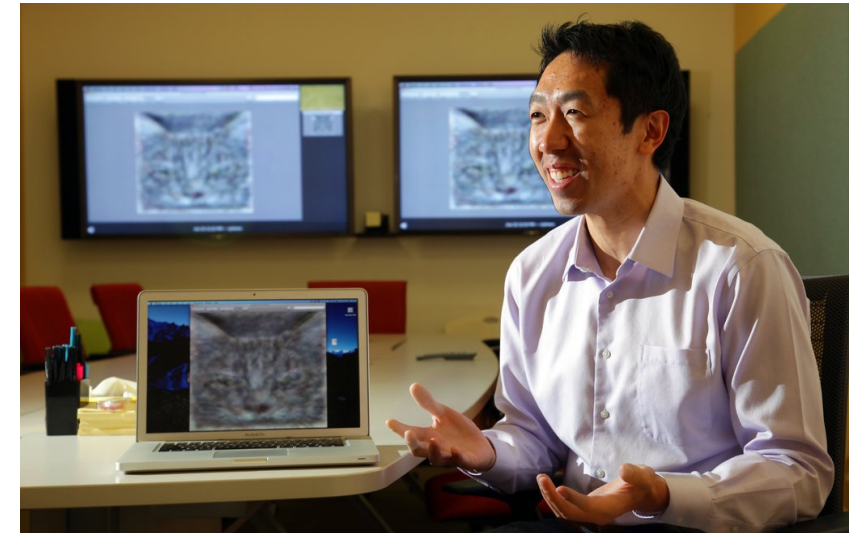
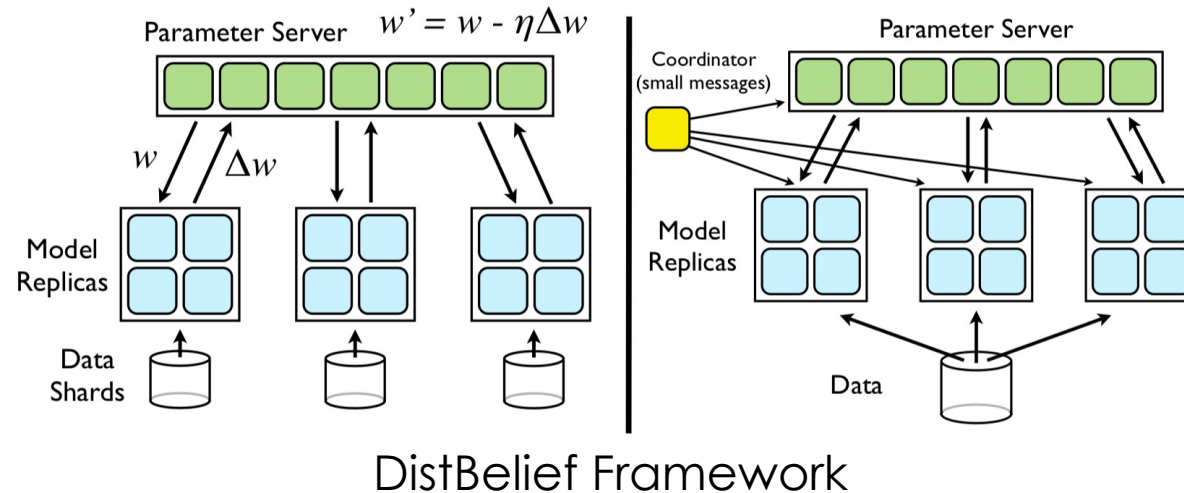
- **Context:** A **vision paper** that outlines the **current state of automatic differentiation** techniques and proposes a new functional, typed **intermediate representation** (IR)
- **Key Idea:** Observe convergence of imperative and declarative approaches and draws connections to compilers → argues for the need for a common IR like those found in modern compilers.
- **Contribution:** Frames problem space and range of techniques.
- **Rational for Reading:** condensed context and some insights for future research directions

TensorFlow: A System for Large-Scale Machine Learning

Large fraction of Google Brain team under Jeff Dean

Context

- Need for distributed training for Deep Learning
- Parameter server abstractions were too general
 - Difficult to use

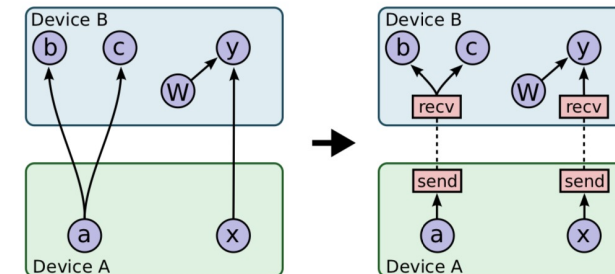


- Theano not designed for distributed setting

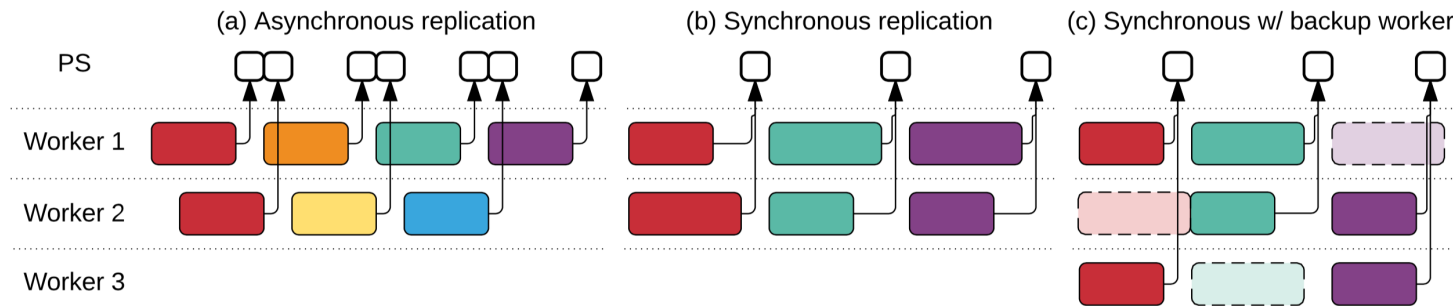
Big Ideas

- Adopts a **dataflow programming** abstraction
 - Inspired by distributed **data processing systems** (@ google)
 - Resulting abstraction is very **similar to Theano**
- Fine grained placement of operations on devices

```
c = []
for d in ['/device:GPU:2', '/device:GPU:3']:
    with tf.device(d):
        a = tf.constant([1.0, 2.0, 3.0, 4.0, 5.0, 6.0], shape=[2, 3])
        b = tf.constant([1.0, 2.0, 3.0, 4.0, 5.0, 6.0], shape=[3, 2])
        c.append(tf.matmul(a, b))
```



- Support multiple distributed concurrency protocols



Recent advances in TensorFlow

- **Keras** : high-level layer composition API

```
# Define the model sequentially
model = tf.keras.Sequential([
    # Adds two densely-connected layers with 64 units:
    layers.Dense(64, activation='relu', input_shape=(32,)),
    layers.Dense(64, activation='relu'),
    # Add a softmax layer with 10 output units:
    layers.Dense(10, activation='softmax')])

# Setup the model and training routines
model.compile(optimizer=tf.train.AdamOptimizer(0.001),
              loss='categorical_crossentropy',
              metrics=['accuracy'])

# Invent some Data
data = np.random.random((1000, 32))
labels = random_one_hot_labels((1000, 10))
# Train the model
model.fit(data, labels, epochs=10, batch_size=32)
# Make predictions
result = model.predict(data, batch_size=32)
```

What to think about when reading

- Relationship and comparisons to Theano?
- Support for distributed computing and exposed abstraction?
- What are the implications of design decisions on an Eager Execution

Additional Reading

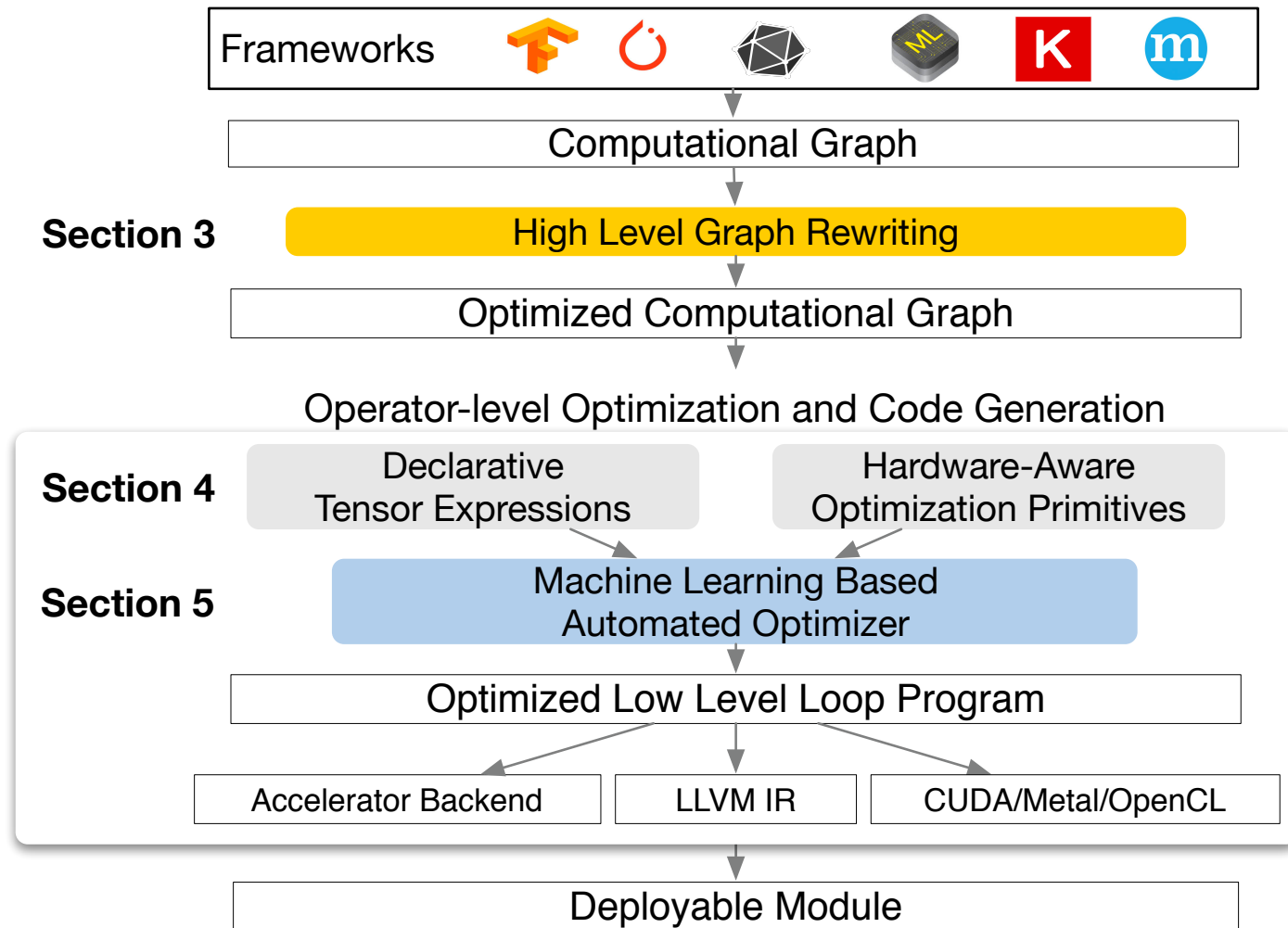
- [TensorFlow: Large-Scale Machine Learning on Heterogeneous Distributed Systems](#)

TVM

Tianqi et. al [OSDI'18]

*Currently visiting in the RISE Lab

TVM



- Originally derived from Halide
 - Leverages similar IR and separation of algorithm from schedule
- Focused on inference workloads

TVM

- Originally derived from Halide
- Leverages similar IR and separation of algorithm from schedule

```
import tvm
```

```
m, n, h = tvm.var('m'), tvm.var('n'), tvm.var('h')  
A = tvm.placeholder((m, h), name='A')  
B = tvm.placeholder((n, h), name='B')
```

Inputs

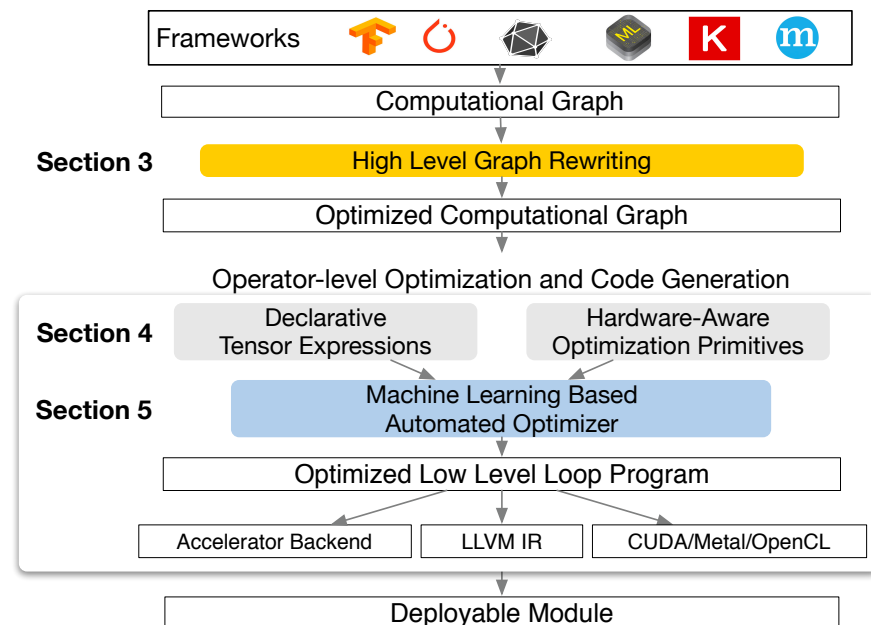
```
k = tvm.reduce_axis((0, h), name='k')  
C = tvm.compute((m, n), lambda i, j: tvm.sum(A[i, k] * B[j, k], axis=k))
```

Shape of C

Computation Rule

```
out = tvm.compute((c, h, w),  
    lambda i, x, y: tvm.sum(data[kc,x+kx,y+ky] * w[i,kx,ky], [kx,ky,kc]))
```

Guess what this describes?



TVM

- Enables declaring new hardware intrinsics
- Simplifies adding support for new hardware

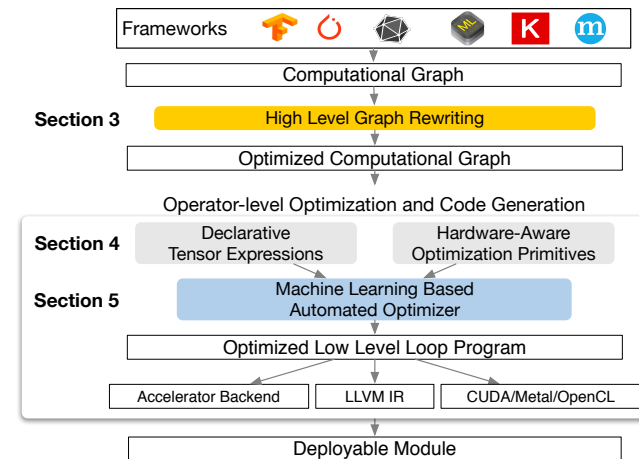
```
w, x = t.placeholder((8, 8)), t.placeholder((8, 8))
k = t.reduce_axis((0, 8))
y = t.compute((8, 8), lambda i, j:
               t.sum(w[i, k] * x[j, k], axis=k))
```

declare behavior

```
def gemm_intrin_lower(inputs, outputs):
    ww_ptr = inputs[0].access_ptr("r")
    xx_ptr = inputs[1].access_ptr("r")
    zz_ptr = outputs[0].access_ptr("w")
    compute = t.hardware_intrin("gemm8x8", ww_ptr, xx_ptr, zz_ptr)
    reset = t.hardware_intrin("fill_zero", zz_ptr)
    update = t.hardware_intrin("fuse_gemm8x8_add", ww_ptr, xx_ptr, zz_ptr)
    return compute, reset, update
```

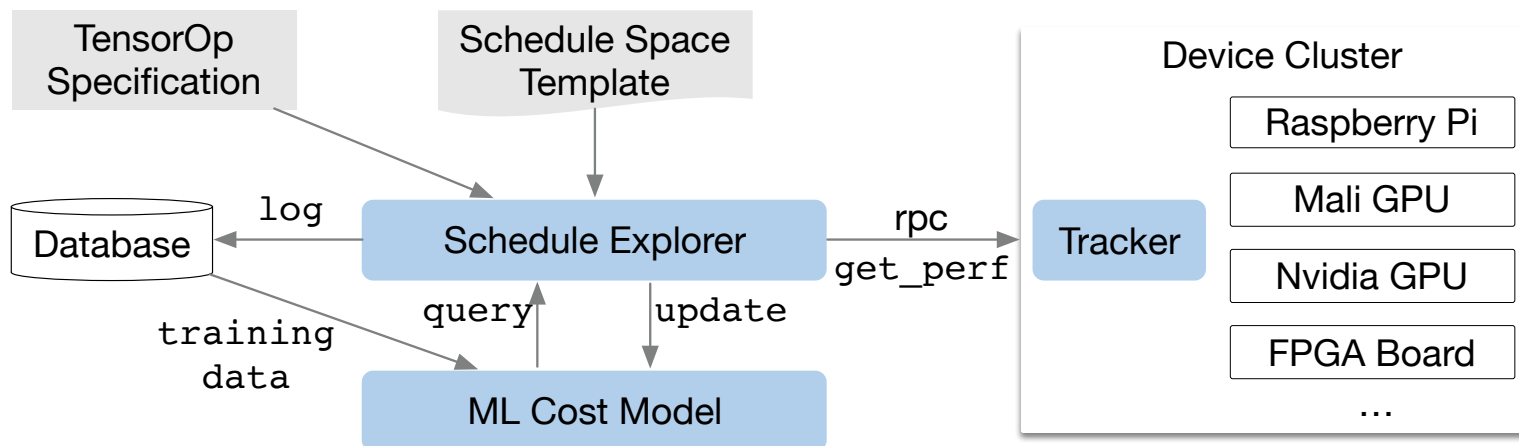
lowering rule to generate
hardware intrinsics to carry
out the computation

```
gemm8x8 = t.decl_tensor_intrin(y.op, gemm_intrin_lower)
```



TVM

➤ Learning based auto-tuner



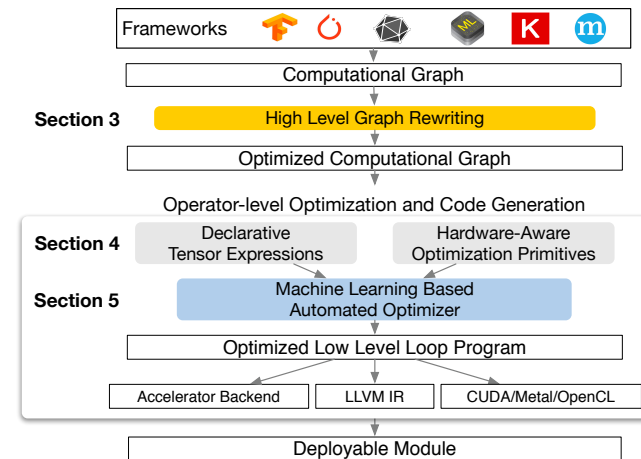
```

for y in range(8):
    for x in range(8):
        C[y][x]=0
        for k in range(8):
            C[y][x]+=A[k][y]*B[k][x]
    
```

(a) Low level AST

	touched memory			outer loop length	
	C	A	B		
y	64	64	64	y	1
x	8	8	64	x	8
k	1	8	8	k	64

(b) Loop context vectors



- Parametrized the AST
- Use Gradient Boosted Trees (GBT) to optimize a “rank loss” to predict the relative order of program runtime

Done!