# Building distributed machine learning pipeline
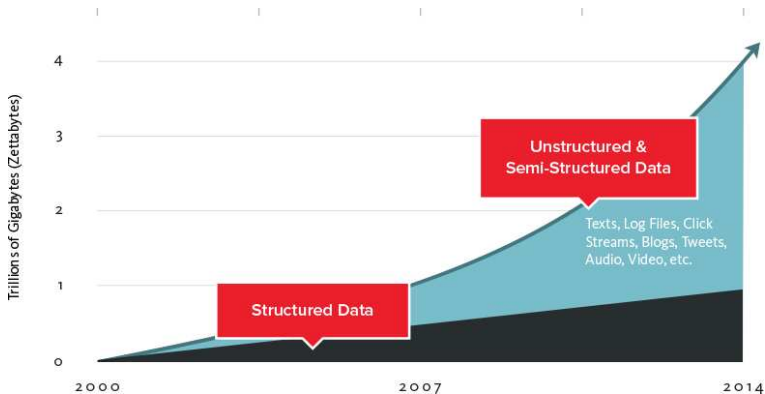
Vojtěch Juránek

JBoss - a division by Red Hat

27. 1. 2017, DEVCONF.CZ, Brno

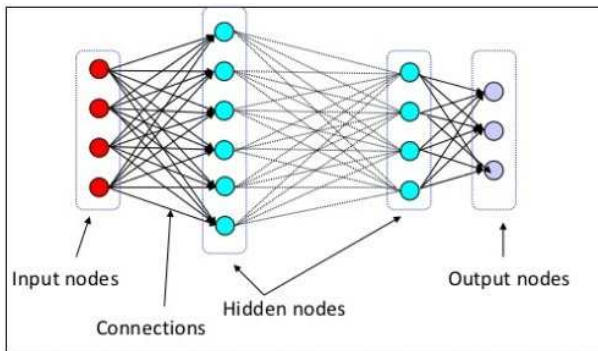# Data today



Source: http://www.couchbase.com/nosql-resources/what-is-no-sql

# Dealing with unstructured data

- There are classes of problems where simple (e.g. linear) models work fine, but there are also classes there liner models fail.

# Deep learning

# Deep learning



Source: http://www.kdnuggets.com/2016/10/deep-learning-key-terms-explained.html

# Big Data → Fast Data

- Pressure to react/process data immediately as it arrives, provide user immediate feedback.
- Big Data → Fast Data

# Big Data → Fast Data

- Pressure to react/process data immediately as it arrives, provide user immediate feedback.
- Big Data → Fast Data

# Big Data → Fast Data

- Pressure to react/process data immediately as it arrives, provide user immediate feedback.
- Big Data → Fast Data

- Process data once it arrives (data/event streaming).
- Do the data analysis with frameworks which keep the data in memory during processing.
- If possible, keep data in memory during whole application stack.
- More in my DevConf 2016 talk (slides).

# Big Data → Fast Data

- Pressure to react/process data immediately as it arrives, provide user immediate feedback.
- Big Data → Fast Data

- Process data once it arrives (data/event streaming).
- Do the data analysis with frameworks which keep the data in memory during processing.
- If possible, keep data in memory during whole application stack.
- More in my DevConf 2016 talk (slides).

# Big Data → Fast Data

- Pressure to react/process data immediately as it arrives, provide user immediate feedback.
- Big Data → Fast Data

- Process data once it arrives (data/event streaming).
- Do the data analysis with frameworks which keep the data in memory during processing.
- If possible, keep data in memory during whole application stack.
- More in my DevConf 2016 talk (slides).

# Big Data → Fast Data

- Pressure to react/process data immediately as it arrives, provide user immediate feedback.
- Big Data → Fast Data

- Process data once it arrives (data/event streaming).
- Do the data analysis with frameworks which keep the data in memory during processing.
- If possible, keep data in memory during whole application stack.
- More in my DevConf 2016 talk (slides).

# What we have so far?

**Problem**

# What we have so far?

## **Problem**

- Data which requires (complicated) non-liner model to sort out.
- We'd like to process and pass incoming data through whole application stack immediately.

# Building neural network pipeline in era of fast data

# Building neural network pipeline in era of fast data

# Deep NN/learning frameworks

- There are many ...
- ... Caffe, CNTK, Deeplearning4j, TensorFlow, ...
- Brief comparison on Wikipedia:
  `https://en.wikipedia.org/wiki/Comparison_of_deep_learning_software`

# Deep NN/learning frameworks

- There are many . . .
- . . . Caffe, CNTK, Deeplearning4j, TensorFlow, . . .
- Brief comparison on Wikipedia:
  `https://en.wikipedia.org/wiki/Comparison_of_deep_learning_software`

# Deep NN/learning frameworks

- There are many . . .
- . . . Caffe, CNTK, Deeplearning4j, TensorFlow, . . .
- Brief comparison on Wikipedia:
  **https://en.wikipedia.org/wiki/Comparison_of_deep_learning_software**

TensorFlow

- Open source software library for machine learning developed by Google Brain team, white paper: arXiv:1603.04467
- Google's second generation machine learning system (first one was DistBelief), open sourced more then year ago (Nov. 2015).
- Used by Google speech recognition, Google photos and other Google products.
- Used also by many other projects, e.g. Mozilla DeepSpeech project (there's a talk by Tilman Kamp about this project here at DevConf on Sunday)

TensorFlow

- Open source software library for machine learning developed by Google Brain team, white paper: arXiv:1603.04467
- Google's second generation machine learning system (first one was DistBelief), open sourced more then year ago (Nov. 2015).
- Used by Google speech recognition, Google photos and other Google products.
- Used also by many other projects, e.g. Mozilla DeepSpeech project (there's a talk by Tilman Kamp about this project here at DevConf on Sunday)
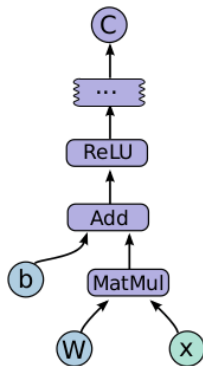
**TensorFlow**

- Open source software library for machine learning developed by Google Brain team, white paper: arXiv:1603.04467
- Google's second generation machine learning system (first one was DistBelief), open sourced more then year ago (Nov. 2015).
- Used by Google speech recognition, Google photos and other Google products.
- Used also by many other projects, e.g. Mozilla DeepSpeech project (there's a talk by Tilman Kamp about this project here at DevConf on Sunday)

TensorFlow

- Open source software library for machine learning developed by Google Brain team, white paper: arXiv:1603.04467
- Google's second generation machine learning system (first one was DistBelief), open sourced more then year ago (Nov. 2015).
- Used by Google speech recognition, Google photos and other Google products.
- Used also by many other projects, e.g. Mozilla DeepSpeech project (there's a talk by Tilman Kamp about this project here at DevConf on Sunday)
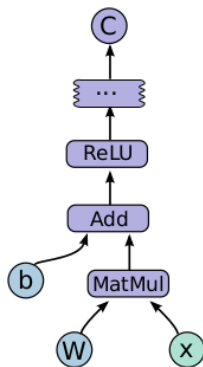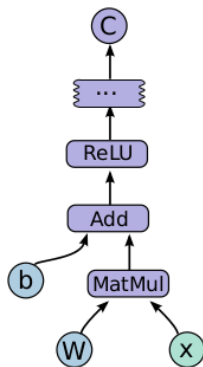
# TensorFlow computation graphs

- **Graph** represents TF computation.
- **Graph nodes** act as mathematical operations.
- **Graph edges** represent matrices (tensors).
- **Session** represents a client accessing particular TF runtime.
- **Variable** is a variable with pre-defined value (pre-defined parameter).
- **Placeholder** is a variable placeholder and its value will be injected into the session.
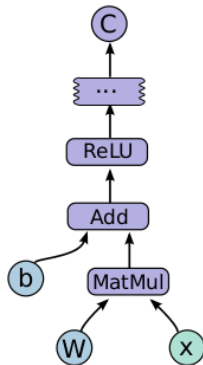- **Checkpoint** use for storing variable/trained model/etc.

# TensorFlow computation graphs

- **Graph** represents TF computation.
- **Graph nodes** act as mathematical operations.
- **Graph edges** represent matrices (tensors).
- **Session** represents a client accessing particular TF runtime.
- **Variable** is a variable with pre-defined value (pre-defined parameter).
- **Placeholder** is a variable placeholder and its value will be injected into the session.
- **Checkpoint** use for storing variable/trained model/etc.

# TensorFlow computation graphs

- **Graph** represents TF computation.
- **Graph nodes** act as mathematical operations.
- **Graph edges** represent matrices (tensors).
- **Session** represents a client accessing particular TF runtime.
- **Variable** is a variable with pre-defined value (pre-defined parameter).
- **Placeholder** is a variable placeholder and its value will be injected into the session.
- **Checkpoint** use for storing variable/trained model/etc.
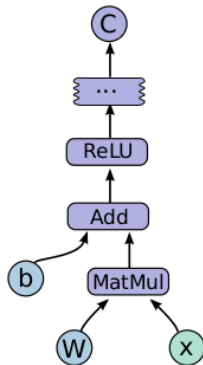
# TensorFlow computation graphs

- **Graph** represents TF computation.
- **Graph nodes** act as mathematical operations.
- **Graph edges** represent matrices (tensors).
- **Session** represents a client accessing particular TF runtime.
- **Variable** is a variable with pre-defined value (pre-defined parameter).
- **Placeholder** is a variable placeholder and its value will be injected into the session.
- **Checkpoint** use for storing variable/trained model/etc.
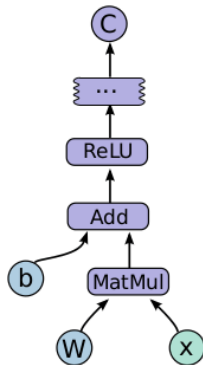
# TensorFlow computation graphs

- **Graph** represents TF computation.
- **Graph nodes** act as mathematical operations.
- **Graph edges** represent matrices (tensors).
- **Session** represents a client accessing particular TF runtime.
- **Variable** is a variable with pre-defined value (pre-defined parameter).
- **Placeholder** is a variable placeholder and its value will be injected into the session.
- **Checkpoint** use for storing variable/trained model/etc.
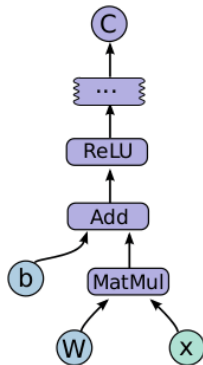
# TensorFlow computation graphs

- **Graph** represents TF computation.
- **Graph nodes** act as mathematical operations.
- **Graph edges** represent matrices (tensors).
- **Session** represents a client accessing particular TF runtime.
- **Variable** is a variable with pre-defined value (pre-defined parameter).
- **Placeholder** is a variable placeholder and its value will be injected into the session.
- **Checkpoint** use for storing variable/trained model/etc.
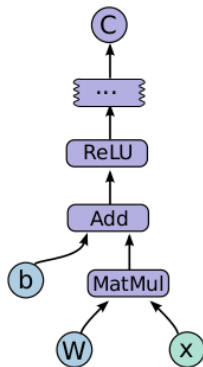
# TensorFlow computation graphs

- **Graph** represents TF computation.
- **Graph nodes** act as mathematical operations.
- **Graph edges** represent matrices (tensors).
- **Session** represents a client accessing particular TF runtime.
- **Variable** is a variable with pre-defined value (pre-defined parameter).
- **Placeholder** is a variable placeholder and its value will be injected into the session.
- **Checkpoint** use for storing variable/trained model/etc.

# TensorFlow computation graphs

```
1  import tensorflow as tf
2
3  b = tf.Variable(tf.zeros([100]))
4  W = tf.Variable(tf.random_uniform \
5    ([784,100],-1,1))
6  x = tf.placeholder(name="x")
7  relu = tf.nn.relu(tf.matmul(W, x) + b)
8  C = [...]
9
10 s = tf.Session()
11 for step in xrange(0, 10):
12   x_in = [...] #construct 100-D input array
13   result = s.run(C, feed_dict={x: x_in})
14   print step, result
```

# Brief list of some others TF features

- Storing graphs and checkpoints into language neutral files (`protobuf` format).
- Supports computation on CPU as well as on GPU, optionally supports CUDA.

```
with tf.Session() as sess:
  with tf.device("/gpu:1"):
```

- Can be run across the cluster, uses `grpc` for communication

```
with tf.device("/job:ps/task:0"):
  weights_1 = tf.Variable(...)
  biases_1 = tf.Variable(...)
[...]
with tf.device("/job:worker/task:7"):
  layer_1 = tf.nn.relu(tf.matmul(input, weights_1) +
      biases_1)
  logits = tf.nn.relu(tf.matmul(layer_1, weights_2) +
      biases_2)
[...]
with tf.Session("grpc://worker7.example.com:2222") as sess:
  sess.run(train_op)
```

# Brief list of some others TF features

- Storing graphs and checkpoints into language neutral files (`protobuf` format).
- Supports computation on CPU as well as on GPU, optionally supports CUDA.

```
1 with tf.Session() as sess:
2   with tf.device("/gpu:1"):
```

- Can be run across the cluster, uses `grpc` for communication

```
1 with tf.device("/job:ps/task:0"):
2   weights_1 = tf.Variable(...)
3   biases_1 = tf.Variable(...)
4 [...]
5 with tf.device("/job:worker/task:7"):
6   layer_1 = tf.nn.relu(tf.matmul(input, weights_1) +
        biases_1)
7   logits = tf.nn.relu(tf.matmul(layer_1, weights_2) +
        biases_2)
8 [...]
9 with tf.Session("grpc://worker7.example.com:2222") as sess:
10   sess.run(train_op)
```

# Brief list of some others TF features

- Storing graphs and checkpoints into language neutral files (`protobuf` format).
- Supports computation on CPU as well as on GPU, optionally supports CUDA.

```
1 with tf.Session() as sess:
2   with tf.device("/gpu:1"):
```

- Can be run across the cluster, uses `grpc` for communication

```
1 with tf.device("/job:ps/task:0"):
2   weights_1 = tf.Variable(...)
3   biases_1 = tf.Variable(...)
4 [...]
5 with tf.device("/job:worker/task:7"):
6   layer_1 = tf.nn.relu(tf.matmul(input, weights_1) +
        biases_1)
7   logits = tf.nn.relu(tf.matmul(layer_1, weights_2) +
        biases_2)
8 [...]
9 with tf.Session("grpc://worker7.example.com:2222") as sess:
10   sess.run(train_op)
```

# TensorFlow integration with Infinispan/Java

- Typically you create the model in Python and run it in C++ or Java app.
- TensorFlow Java API support still TBD, see TF issue #5 and issue #3
- C++ API via JNI can be used as a workaround for Java.
- You can use some existing library, e.g. javacpp-presets.

```
1  GraphDef graph = new GraphDef();
2  ReadBinaryProto(Env.Default(), modelPath, graph);
3  SessionOptions options = new SessionOptions();
4  this.session = new Session(options);
5  Status status = session.Create(graph);
6  [...]
7  Tensor img = new Tensor(DT_FLOAT, new TensorShape(1, image.
       length));
8  FloatBuffer imgBuff = img.createBuffer();
9  imgBuff.put(image);
10 [...]
11 TensorVector outputs = new TensorVector();
12 Status status = session.Run(new StringTensorPairVector(new
       String[] { "images" }, new Tensor[] { img }),
13 new StringVector("softmax_linear/logits"), new StringVector("
       softmax_linear/logits"), outputs);
```

# In-memory data grid: Infinispan

**http://infinispan.org/**

- Data grid platform, written in Java
- In-memory No-SQL key-value data store, (optionally) schema-less
- Distributed cache - offers massive memory
- Elastic and scalable - can run on hundreds of nodes
- Highly available - no SPOF, resilient to node failures
- Transactional
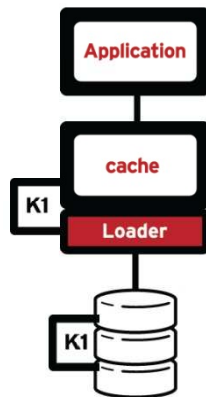- Supports indexing and searching
- Many other features

# Infinispan eviction and cache stores

**Eviction:** removing entries from the cache.

```
ConfigurationBuilder().eviction().size(5).
    strategy(EvictionStrategy.LRU)
```

**Cache store:** a way how to store cache content in some external (persistent) storage.
There are many, JDBC, JPA, clouds, LevelDB, Cassandra . . . and also **Ceph** cache store.
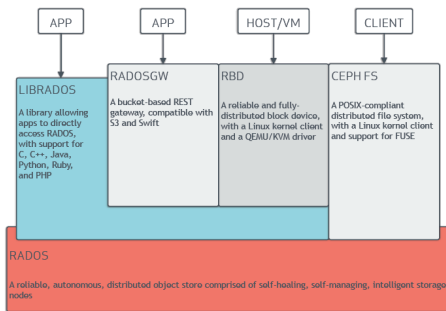
# Ceph



**http://ceph.com/**

- Distributed object storage.
- Provides interfaces to object, block and file system storage in unified data cluster.
- Scalable to the exabyte level.
- Highly available - no SPOF, resilient to node failures.
- Open source.

# Ceph architecture and Infinispan cache store

- `rados` is distributed object store
- Infinispan cache store accesses `rados` directly via `librados` (using Java JNI client).



Ceph cache store configuration:
(details on https://github.com/vjuranek/infinispan-cachestore-ceph).

```
1  <local-cache name="cachestore" start="EAGER">
2      <store class="org.infinispan.persistence.ceph.CephStore">
3          <property name="monitorHost">192.168.122.145:6789</
               property>
4          <property name="userName">admin</property>
5          <property name="key">mykey</property>
6      </store>
7  </local-cache>
```

Sources available on **https://github.com/vjuranek/tf-ispn-demo**

## **NN "Hello world" - MNIST data sample**

- Recognition of hand written digits.
- Besides simplicity, TF has very detailed tutorial for MNIST (for beginners as well as for experts)

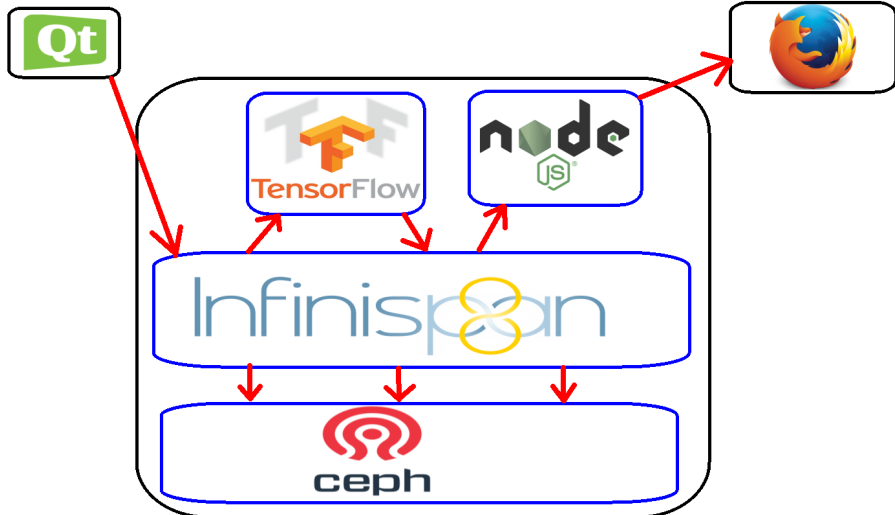Sources available on **https://github.com/vjuranek/tf-ispn-demo**

## **NN "Hello world" - MNIST data sample**

- Recognition of hand written digits.
- Besides simplicity, TF has very detailed tutorial for MNIST (for beginners as well as for experts)

Sources available on **https://github.com/vjuranek/tf-ispn-demo**

## **NN "Hello world" - MNIST data sample**

- Recognition of hand written digits.
- Besides simplicity, TF has very detailed tutorial for MNIST (for beginners as well as for experts)

# "Hello world" Demo

# Summary

- Building pipeline for complex data processing in real time can be quite easy if you use the right tools.
- TensorFlow is very powerful machine learning framework.
- Infinispan is real middleware which can glue together various pieces of your application stack and server as its backbone.

# Summary

- Building pipeline for complex data processing in real time can be quite easy if you use the right tools.
- TensorFlow is very powerful machine learning framework.
- Infinispan is real middleware which can glue together various pieces of your application stack and server as its backbone.

# Summary

- Building pipeline for complex data processing in real time can be quite easy if you use the right tools.
- TensorFlow is very powerful machine learning framework.
- Infinispan is real middleware which can glue together various pieces of your application stack and server as its backbone.

# Question?

http://infinispan.org/

# Thank you for your attention!