Hands-on Machine Learning with Kafka-based Streaming Pipelines

Strata Data, London, 2019

Boris Lublinsky and Dean Wampler, Lightbend

boris.lublinsky@lightbend.com dean.wampler@lightbend.com



©Copyright 2017-2019, Lightbend, Inc.

Apache 2.0 License. Please use as you see fit, but attribution is requested.

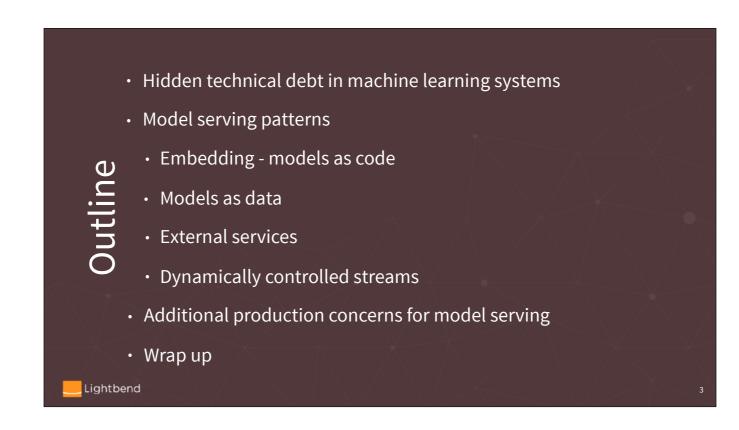
If you have not done so already, download the tutorial from GitHub

https://github.com/lightbend/model-serving-tutorial

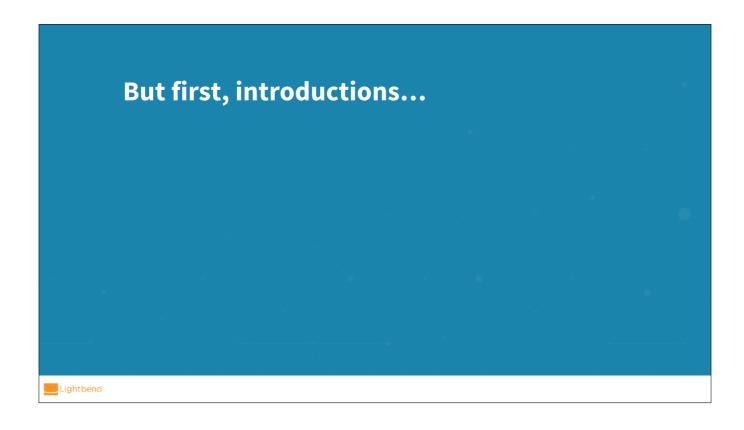
See the README for setup instructions.

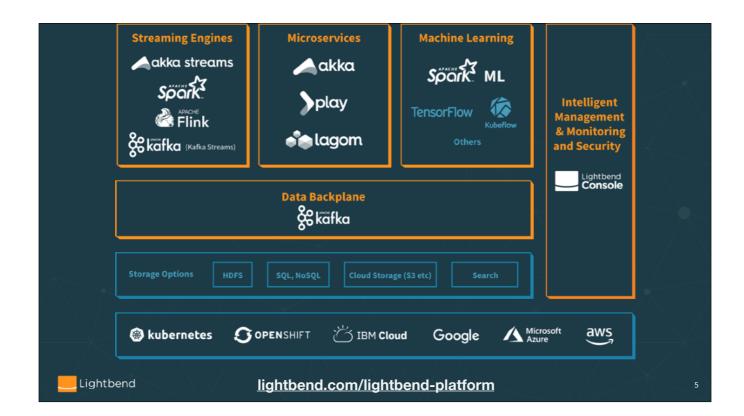
These slides are in the presentation folder.

Lightbend

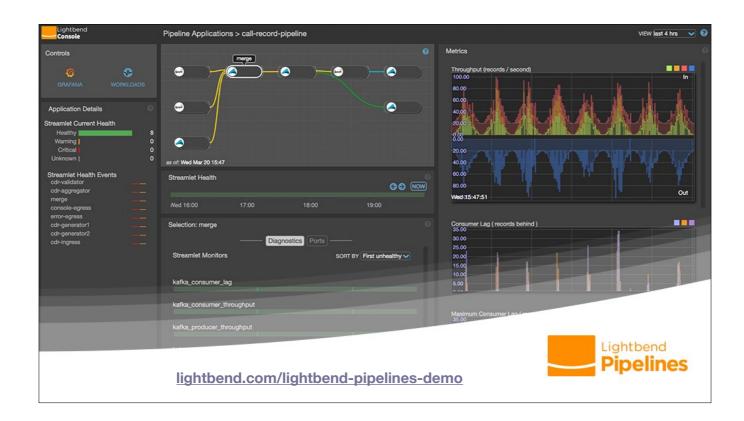


As we'll see dynamically controlled streams will revisit models as data





<marketing>Lightbend has been thinking about the development of microservices and streaming data applications for a while. Lightbend Platform is a commercial distribution of open-source and commercial tools to accelerate development, deployment, and management of these applications.



<marketing>In particular, Lightbend Pipelines is a new capability in Lightbend Platform for writing, deploying, monitoring, and managing streaming data pipelines.
marketing>



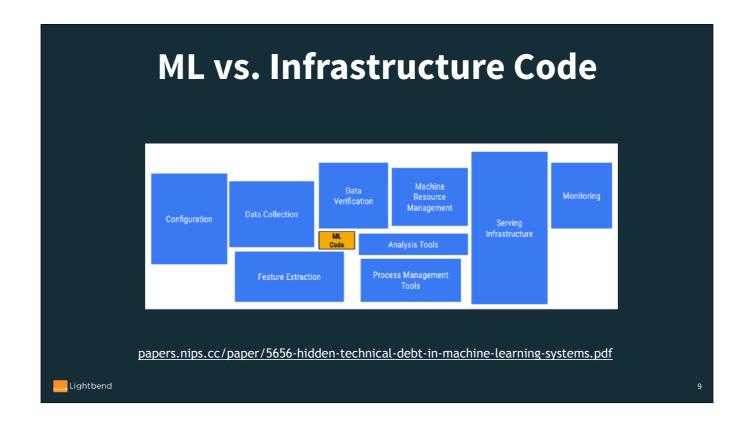
We can only touch on architectures, design principles, and implementation idioms in this tutorial. For more extensive coverage, see these free reports written by Boris and Dean, which are published by O'Reilly.

• Hidden technical debt in machine learning systems

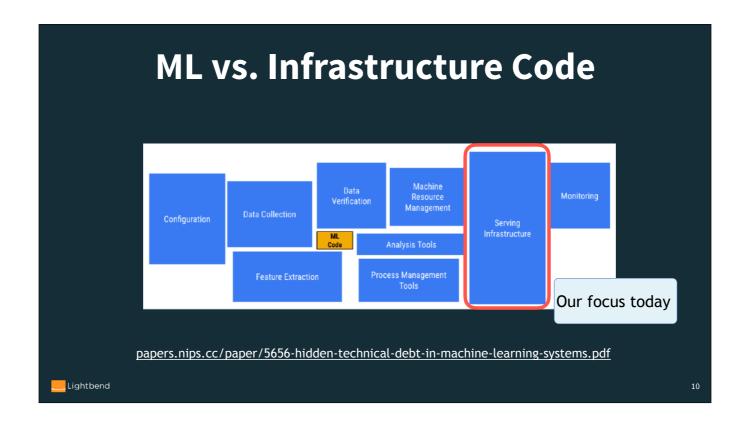
- Model serving patterns
 - Embedding models as code
 - Models as data
 - External services
 - Dynamically controlled streams
- Additional production concerns for model serving
- Wrap up

Lightbend

œ Ì



This paper discusses the challenges of real-world use of ML. The actual ML code is just a small part of the overall infrastructure and capabilities required.



This tutorial discusses a few aspects of that picture, focusing on architectures for *serving* models in productions. The architectures use minimal coupling to tool chains for *training* (e.g., through a Kafka topic), which enables a wide range of *training* options, making it easier to use your favorite data science tool chain for training.

• Hidden technical debt in machine learning systems

- Model serving patterns
 - Embedding models as code
 - Models as data
 - External services
 - Dynamically controlled streams
- Additional production concerns for model serving
- Wrap up

Lightbend

Model Serving Architectures

- Embedding model as *code*, deployed into a stream engine
- Model as *data* easier dynamic updates
- Model Serving as a service use a separate service, access from the streaming engine
- Dynamically controlled streams one way to implement model as data in a streaming engine

Lightbend

10

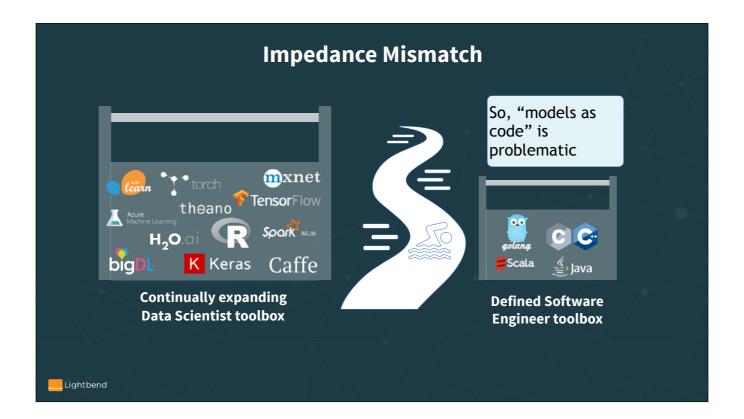
You can embed code for models, but we'll see this has many disadvantages. It's better to treat models as data, which allows you to update it as the "world" changes. We'll focus on two approaches for models as data.

Embedding: Model as Code

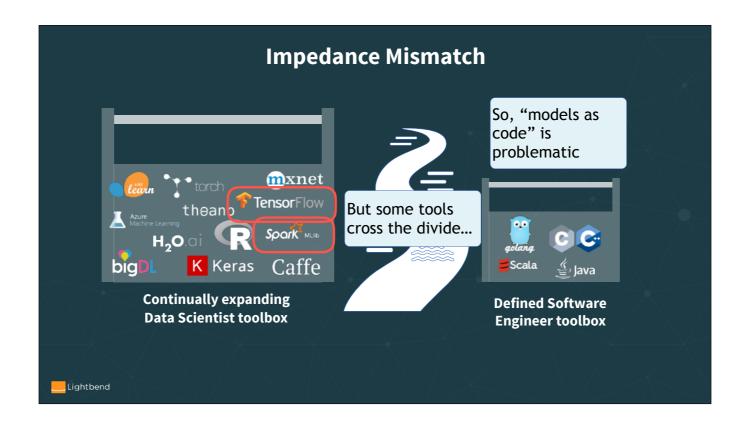
- Implement the model as source code
- The model code is linked into the streaming application at build time

Why is this problematic?

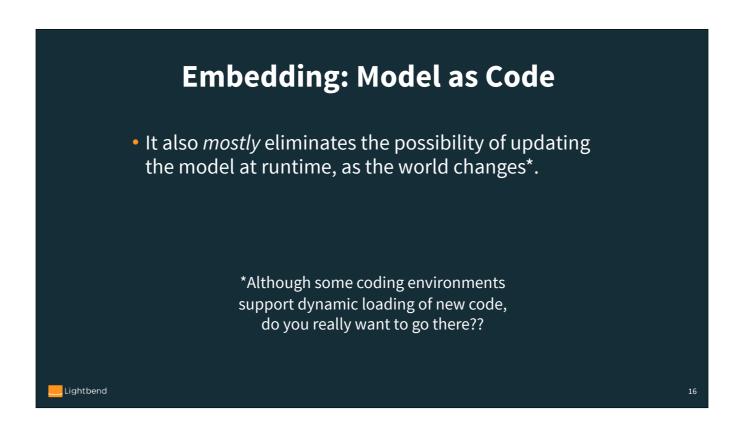
Lightbend



In his talk at the last Flink Forward, Ted Dunning discussed the fact that with multiple tools available to Data scientists, they tend to use different tools for solving different problems and as a result they are not very keen on tools standardization. This creates a problem for software engineers trying to use "proprietary" model serving tools supporting specific machine learning technologies. As data scientists evaluate and introduce new technologies for machine learning, software engineers are forced to introduce new software packages supporting model scoring for these additional technologies.



We loose portability if we use tools crossing divide

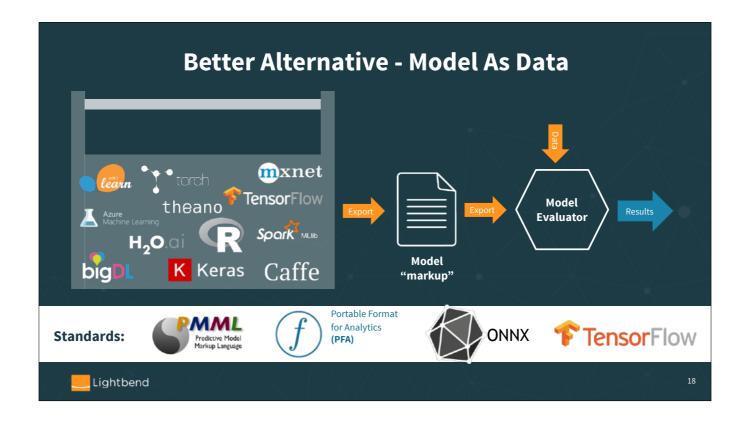


Most dynamically-typed languages, like Python, allow runtime loading of new code. Even the JVM supports this. However, it's much easier to introduce bugs, security holes (SQL injection attacks anyone??), etc. We'll also discuss other production concerns later, like auditing and data governance, which are harder to support using models as code.

• Hidden technical debt in machine learning systems

- Model serving patterns
 - Embedding models as code
 - Models as data
 - External services
 - Dynamically controlled streams
- Additional production concerns for model serving
- Wrap up

Lightbend



In order to overcome these differences, the Data Mining Group has introduced two standards - Predictive Model Markup Language (PMML) and Portable Format for Analytics (PFA), both suited for description of the models that need to be served. Introduction of these models led to creation of several software products dedicated to "generic" model serving, for example Openscoring, Open data group, etc.

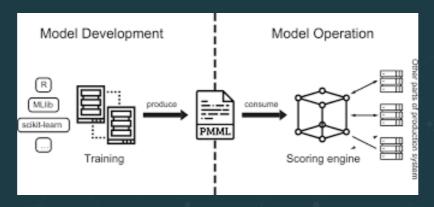
ONNX is a new standard for deep learning models, but it is not supported by TensorFlow.

TensorFlow itself is a de facto standard for ML, because it is so widely used for both training and serving, even though it uses proprietary formats.

The result of this standardization is creation of the open source projects, supporting these formats - JPMML and Hadrian which are gaining more and more adoption for building model serving implementations, for example ING, R implementation, SparkML support, Flink support, etc. TensorFlow also released TensorFlow java APIs, which are used in a Flink TensorFlow

PMML

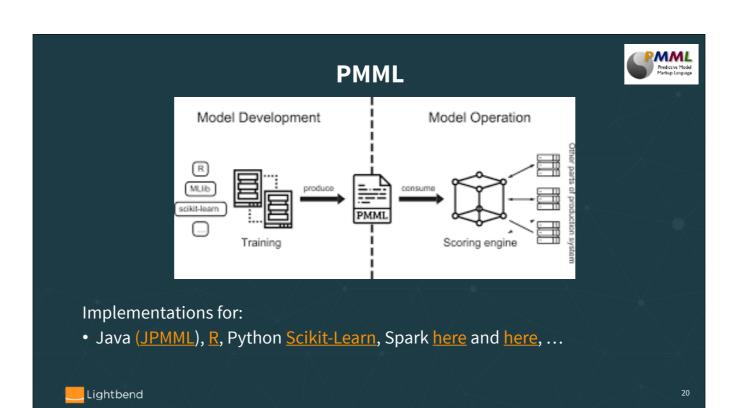




<u>Predictive Model Markup Language</u> (PMML) is an XML-based language that enables the definition and sharing of predictive models between applications.

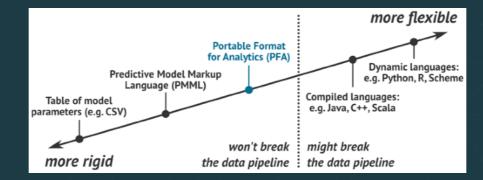
https://www.wismutlabs.com/blog/agile-data-science-2/

Lightbend



PFA





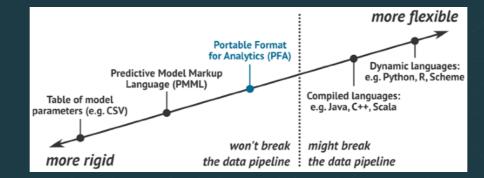
<u>Portable Format for Analytics</u> (PFA) is an emerging standard for statistical models and data transformation engines. PFA combines the ease of portability across systems with algorithmic flexibility: models, pre-processing, and post-processing are all functions that can be arbitrarily composed, chained, or built into complex workflows.



http://dmg.org/pfa/docs/motivation/

PFA





Implementations for:

• Java (<u>Hadrian</u>), R (<u>Aurelius</u>), Python (<u>Titus</u>), Spark (<u>Aardpfark</u>), ...

Lightbend

วว

ONNX





Open Neural Networks Exchange (ONNX) is an open standard format of machine learning models to offer interoperability between various Al frameworks. Led by Facebook, Microsoft, and AWS.

https://azure.microsoft.com/en-us/blog/onnx-runtime-for-inferencing-machine-learning-models-now-in-preview/

Lightbend

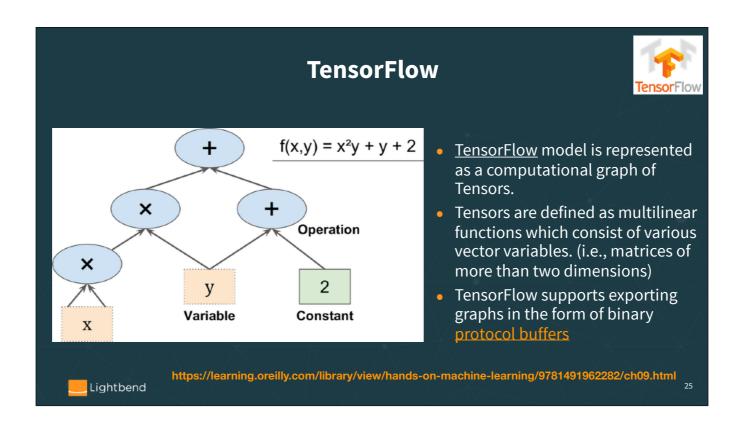
ONNX



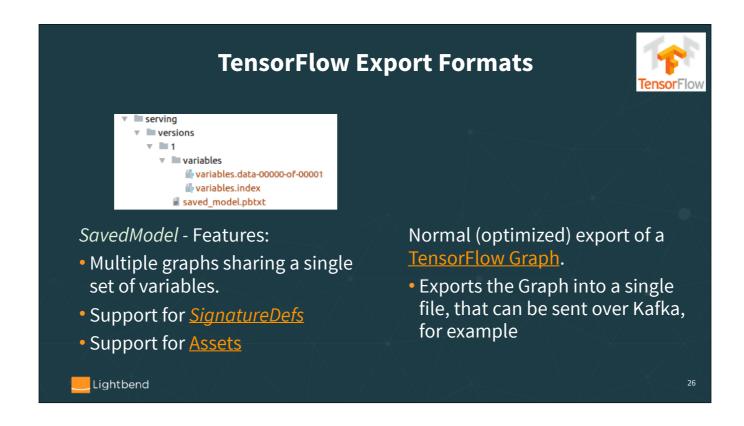


- Supported Tools page.
- Converters for Keras, CoreML, LightGBM, Scikit-Learn,
- PyTorch,
- third-party support for **TensorFlow**

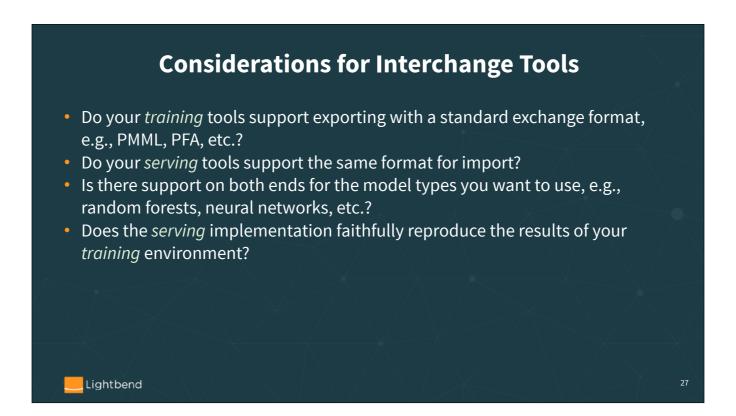
Lightbend



The mathematical concept of a Tensor is more sophisticated than just arrays of more than two dimensions.



Be careful about sending the whole graph into a Kafka topic. Kafka messages are normally not supposed to be large. The upper limit is about 1MB and even that's larger than you should normally use. One alternative, if the model is really big, is to write to a filesystem and send the path in the Kafka message ("pass by reference").

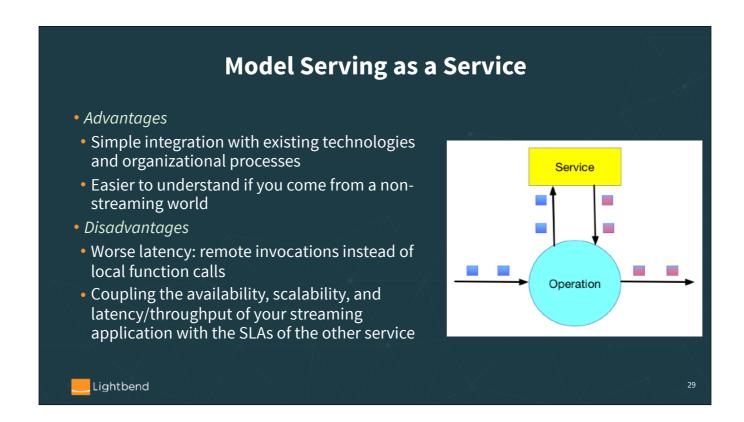


Implementations are spotty, both in terms of import, export coverage, model types, and faithfully reproducing the same results on both ends. I.e., if you use Scikit-Learn to train a random forest, do you get the same results if you run it in SparkML, after exchanging with PMML?

• Hidden technical debt in machine learning systems

- Model serving patterns
 - Embedding models as code
 - Models as data
 - External services
 - Dynamically controlled streams
- Additional production concerns for model serving
- Wrap up

Lightbend



Implementations are spotty, both in terms of import, export coverage, model types, and faithfully reproducing the same results on both ends. I.e., if you use Scikit-Learn to train a random forest, do you get the same results if you run it in SparkML, after exchanging with PMML?

Model Serving as a Service: Challenges

- Launch ML runtime graphs, scale up/down, perform rolling updates
- Infrastructure optimization for ML
- Latency and throughput optimization
- Connect to business apps via various APIs, e.g. REST, gRPC
- Allow Auditing and clear versioning
- Integrate into Continuous Integration (CI)
- Allow Continuous Deployment (CD)
- Provide Monitoring

adapted from https://github.com/SeldonIO/seldon-core/blob/master/docs/challenges.md

Lightbend

30

The SeldonIO project has a nice list of the challenges you face when doing ML Serving as a Service

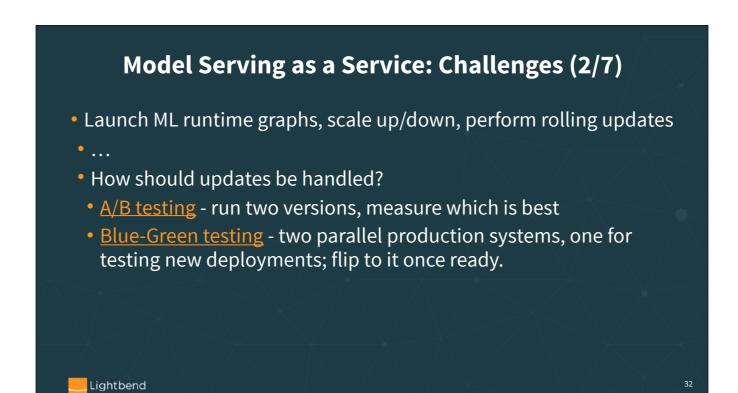
Model Serving as a Service: Challenges (1/7)

- Launch ML runtime graphs, scale up/down, perform rolling updates
- The classic issues for deploying any modern *microservice*.
- What procedures does my organization use to control what's deployed to production and when?
- How do monitor performance?
- When it's necessary to scale up or down, how...?
- How should updates be handled?

Lightbend

3

Let's go through these in a bit more detail.

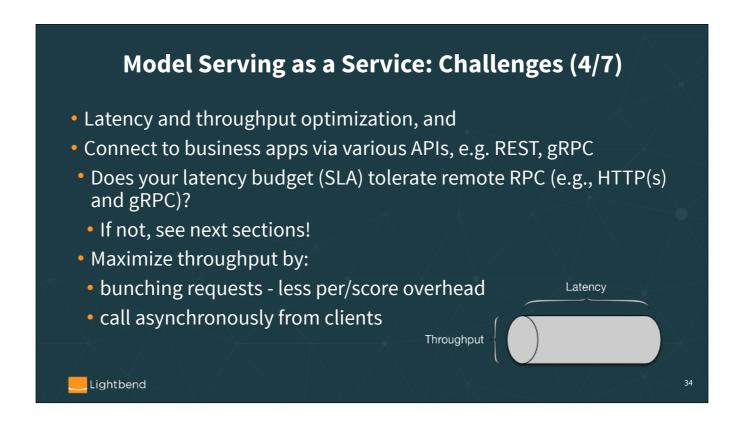


For updates, we'll ignore for now the case where I might run multiple models in production. We'll return to this later.

Model Serving as a Service: Challenges (3/7)

- Infrastructure optimization for ML
- Closer collaboration with Data Science to:
- feed live data for training,
- serve up-to-date models.
- Model training can be very compute-intensitive
- Model serving algorithms can be expensive
- Scoring is less deterministic than other services:
- How do you know when it's wrong?

Lightbend



Discuss two of the bullet points here...

Model Serving as a Service: Challenges (5/7) Allow Auditing and clear versioning When was a particular model used to score a particular record? E.g., "why was my loan application rejected?" Which model was used? When was it deployed? Regulatory compliance

The SeldonIO project has a nice list of the challenges you face when doing ML Serving as a Service

Model Serving as a Service: Challenges (6/7)

- Integrate into Continuous Integration (CI)
- Allow Continuous Deployment (CD)
- Standard and modern techniques for build, test, and deployment
- Probably new to the Data Science teams,
- but hopefully familiar to the Data Engineers!
- This also helps with auditing and versioning requirements

Lightbend

36

Consider two bullet points together, because they are closely related.

Model Serving as a Service: Challenges (7/7)

- Provide Monitoring
- Standard production tool for *observability*
- What's the throughput and latency for scoring?
- Is throughput keeping up with the data volume?
- Is latency meeting our SLAs?
- When will I know that a problem has arisen?
- Alerting essential!
- Practice fault isolation and recovery scenarios

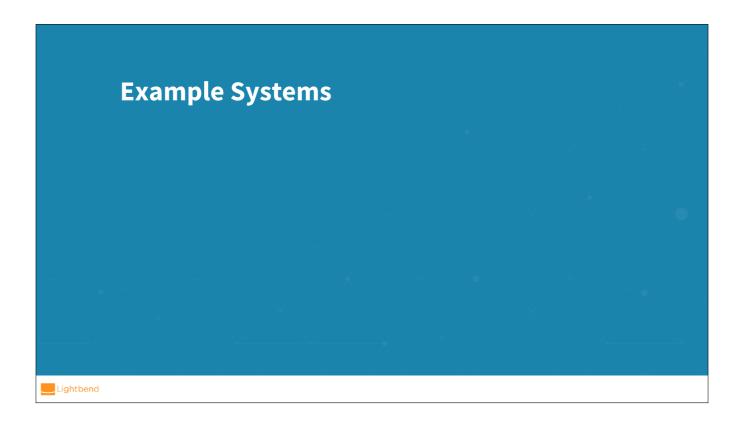
Lightbend

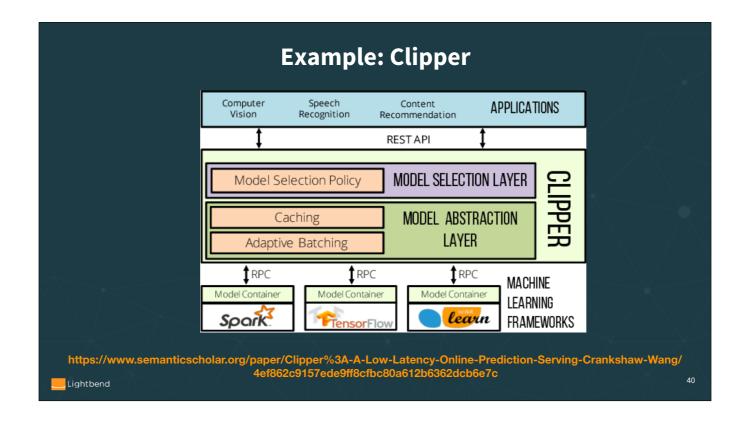
37

Model Serving as a Service: Challenges

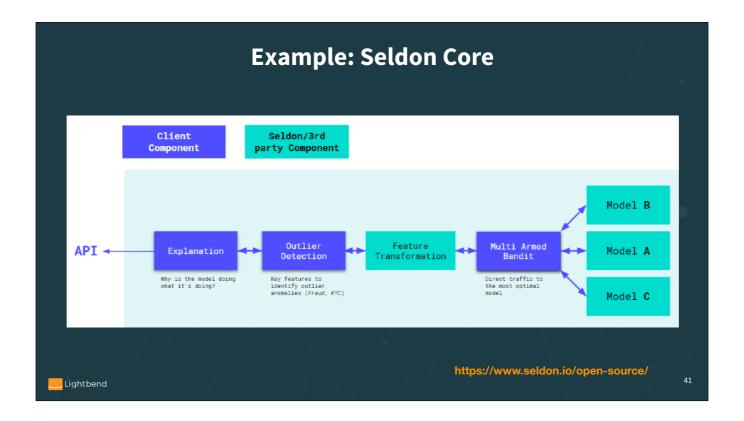
- Launch ML runtime graphs, scale up/down, performeg updates
 Infrastructure optimization for ML
 Latency and throughput optimization
 Connect to business are to larious APIXES RIA; gRPC
 Allow Auditing Mic Continuous Eeg Quon (CI)
 Allow Continuous Eeg Quon (CI)
 Provide O Itoring

Lightbend



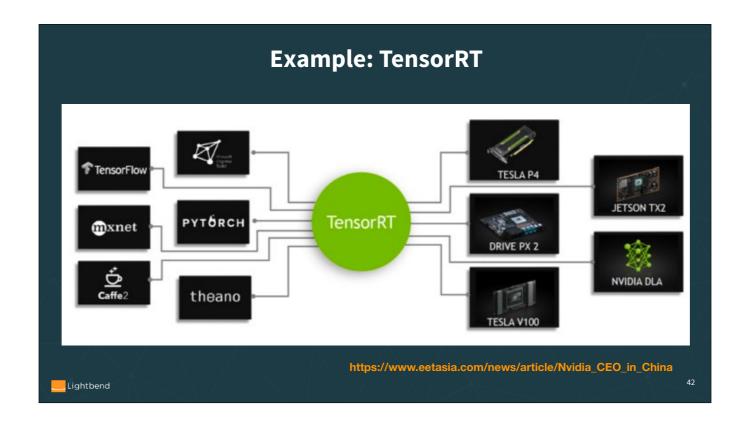


Clipper introduces a modular architecture to simplify model deployment across frameworks and applications. Furthermore, by introducing caching, batching, and adaptive model selection techniques, Clipper reduces prediction latency and improves prediction throughput, accuracy, and robustness without modifying the underlying machine learning frameworks.



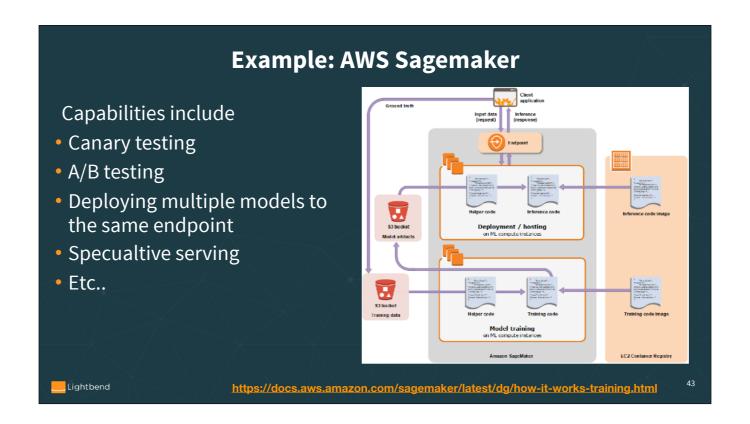
Main characteristics:

- Allow data scientists to create models using any machine learning toolkit or programming language.
- Expose machine learning models via REST and gRPC for easy integration into business apps.
- Allow complex runtime inference graphs to be deployed as microservices. Composed of: models, routers, combiners and transformers.
- Handle full lifecycle management of the deployed model. Updating Scaling Monitoring Security.



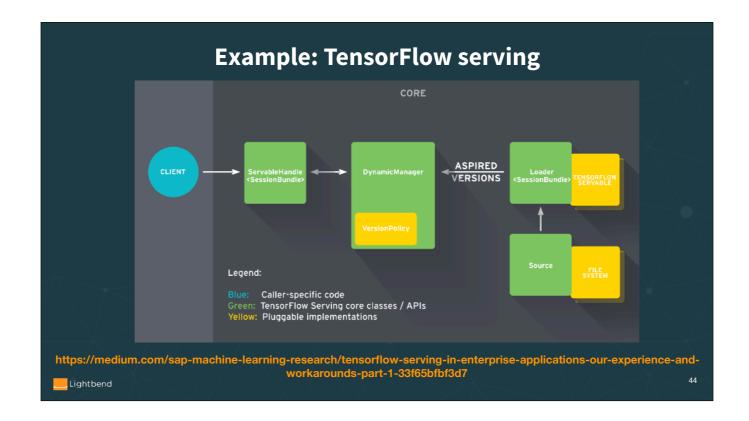
NVIDIA TensorRT™ is a platform for high-performance deep learning inference. It includes a deep learning inference optimizer and runtime that delivers low latency and high-throughput for deep learning inference applications. TensorRT-based applications perform up to 40x faster than CPU-only platforms during inference. With TensorRT, you can optimize neural network models trained in all major frameworks, calibrate for lower precision with high accuracy, and finally deploy to hyperscale data centers, embedded, or automotive product platforms.

TensorRT provides INT8 and FP16 optimizations for production deployments of deep learning inference applications such as video streaming, speech recognition, recommendation and natural language processing. Reduced precision inference significantly reduces application latency, which is a requirement for many real-time services, auto and embedded applications.

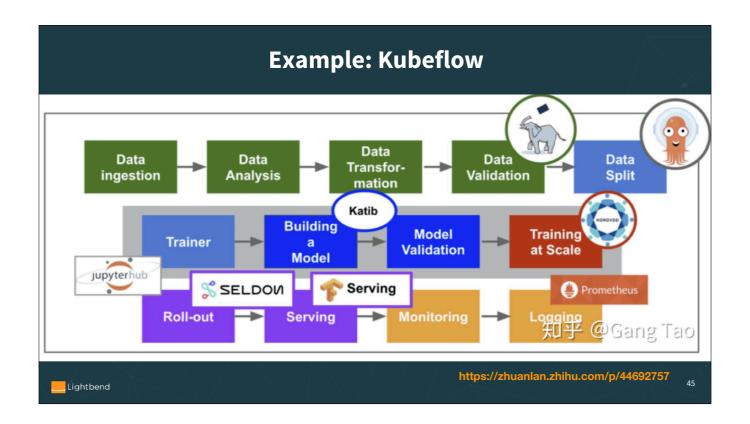


Amazon SageMaker is a fully managed machine learning service. With Amazon SageMaker, data scientists and developers can quickly and easily build and train machine learning models, and then directly deploy them into a production-ready hosted environment. It provides an integrated Jupyter authoring notebook instance for easy access to your data sources for exploration and analysis, so you don't have to manage servers. It also provides common machine learning algorithms that are optimized to run efficiently against extremely large data in a distributed environment. With native support for bring-your-own-algorithms and frameworks, Amazon SageMaker offers flexible distributed training options that adjust to your specific workflows. Deploy a model into a secure and scalable environment by launching it with a single click from the Amazon SageMaker console. Training and hosting are billed by minutes of usage, with no minimum fees and no upfront commitments.

A great article about Sagemaker capabilities and its usage can be found in this article: https://medium.com/@julsimon/mastering-the-mystical-art-of-model-deployment-cocafe011175



TensorFlow Serving is a flexible, high-performance serving system for machine learning models, designed for production environments. TensorFlow Serving makes it easy to deploy new algorithms and experiments, while keeping the same server architecture and APIs. TensorFlow Serving provides out-of-the-box integration with TensorFlow models, but can be easily extended to serve other types of models and data.



The Kubeflow project is dedicated to making deployments of machine learning (ML) workflows on Kubernetes simple, portable and scalable. Our goal is not to recreate other services, but to provide a straightforward way to deploy best-of-breed open-source systems for ML to diverse infrastructures. Anywhere you are running Kubernetes, you should be able to run Kubeflow.

Rendezvous Architecture Pattern

Handle ML logistics in a flexible, responsive, convenient, and realistic way.

- For Data:
- Collect data at scale from many sources.
- Preserve raw data so that potentially valuable features are not lost.
- Make data available to many applications (consumers), both on premise and distributed.



https://mapr.com/ebooks/machine-learning-logistics/ch03.html

Rendezvous Architecture Pattern

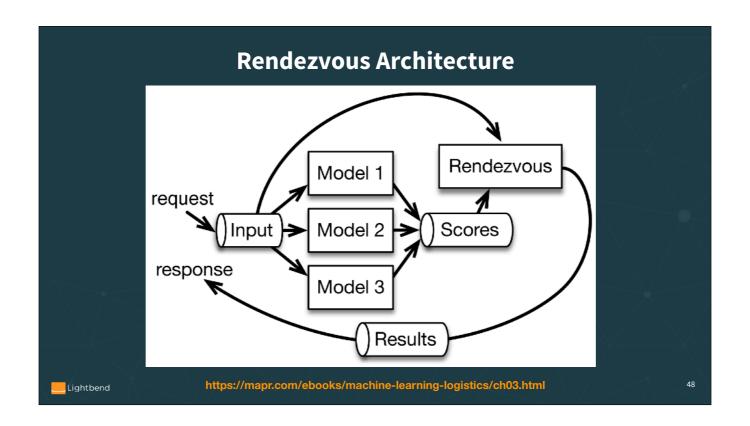
Handle ML logistics in a flexible, responsive, convenient, and realistic way.

- Models:
- Manage multiple models during development and production.
- Improve evaluation methods for comparing models during development and production, including use of reference models for baseline successful performance.
- Have new models poised for rapid deployment.

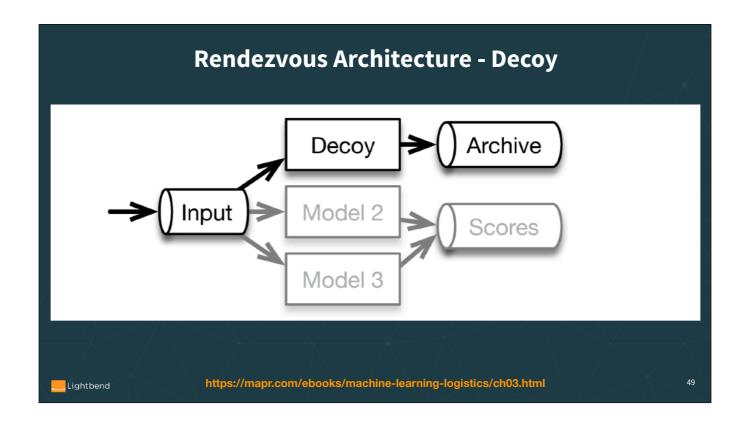
Lightbend

https://mapr.com/ebooks/machine-learning-logistics/ch03.html

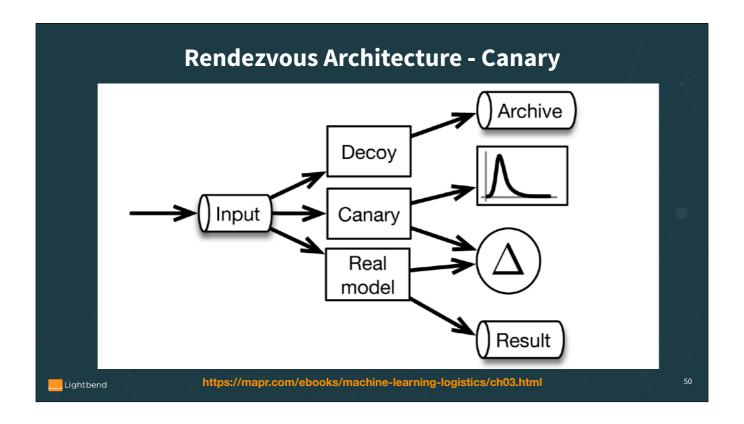
4



All communications here are Kafka or similar log-oriented or message queue alternatives. An input is fanned out to multiple models that serve it simultaneously. The result goes to the scoring function, that combine all of the scoring results (based on certain criteria) and creates combined result.



If we want to collect data for the future training, we can deploy decoy, which is a pass through writing data to the archive. This allows to add data collection without disrupting actual model serving



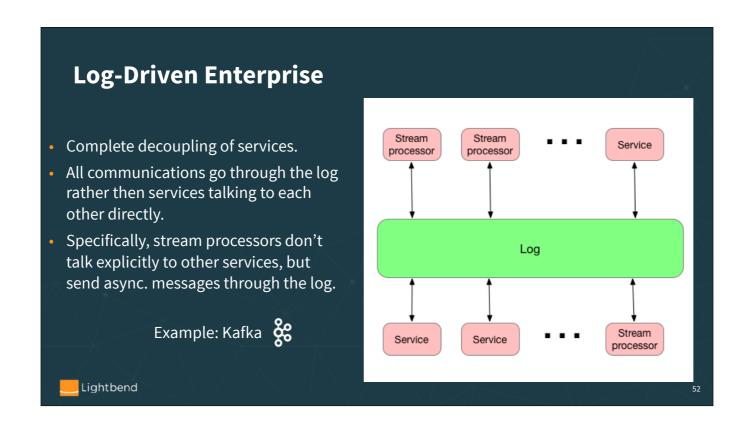
All communications here are Kafka (or its variants). Here we want to validate a new model (canary deployment). This is done deploying a new model to listen to the same topic that the main model and comparing results between new and existing model. If the comparison is good, the new model can be safely deployed.

• Hidden technical debt in machine learning systems

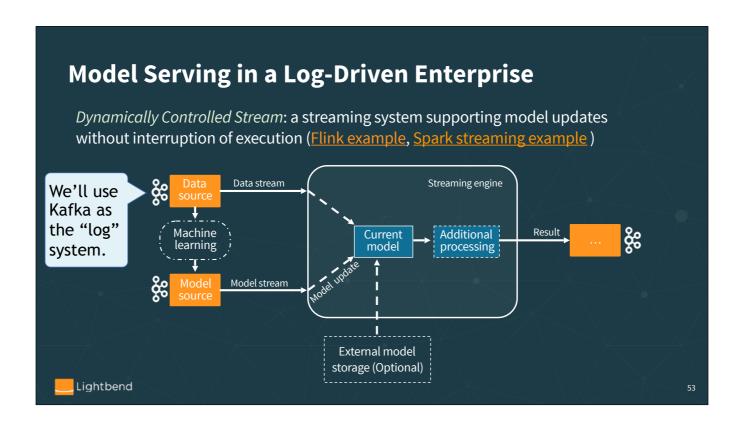
- Model serving patterns
 - Embedding models as code
 - Models as data
 - External services
 - Dynamically controlled streams
- Additional production concerns for model serving
- Wrap up

Lightbend

51



The majority of machine learning implementations are based on running model serving as a Rest service, which might not be appropriate for the high volume data processing or usage of the streaming system, which requires re coding/starting systems for model update, for example, Flink TensorFlow or Flink JPPML.



The majority of machine learning implementations are based on running model serving as a Rest service, which might not be appropriate for the high volume data processing or usage of the streaming system, which requires re coding/starting systems for model update, for example, Flink TensorFlow or Flink JPPML. The name of this pattern was coined by data Artisans (that's the correct spelling...)

```
Model Representation (Protobufs)
// On the wire
                                              See the "protobufs"
syntax = "proto3";
                                              project in the
// Description of the trained model.
                                               example code.
message ModelDescriptor {
 string name = 1;
                       // Model name
 string description = 2; // Human readable
 string dataType = 3;
                      // Data type for which this model is applied.
                      // Model type
 enum ModelType {
                                                   ModelType modeltype = 4;
  TensorFlow = 0;
                                                   oneof MessageContent {
  TensorFlowSAVED = 2;
                                                     // Byte array containing the model
   PMML = 2;
                      // Could add PFA, ONNX, ...
                                                     bytes data = 5;
                                                     string location = 6;
  Lightbend
```

You need a neutral representation format that can be shared between different tools and over the wire. Protobufs (from Google) is one of the popular options. Recall that this is the format used for model export by TensorFlow. Here is an example.

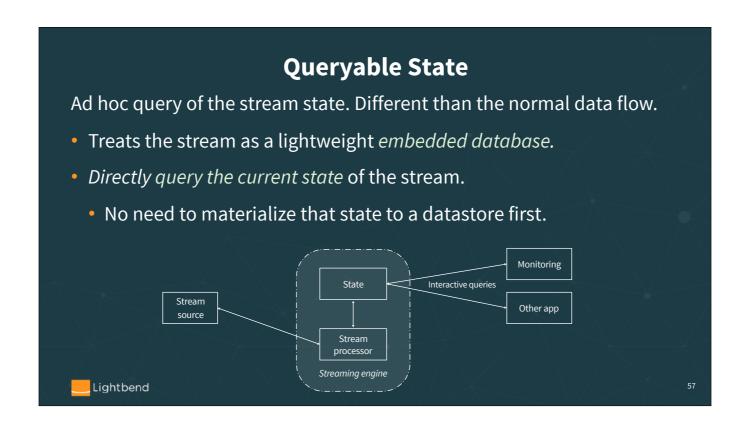
```
Model Code Abstraction (Scala)
                                          [RECORD, RESULT] are
       trait Model[RECORD, RESULT] {
                                          type parameters;
       def score(input : RECORD) : RESULT
                                          compare to Java:
       def cleanup() : Unit
                                          <RECORD, RESULT>
       def toBytes(): Array[Byte]
       def getType : Long
                                                      See the "model" project
                                                      in the example code.
       trait ModelFactory[RECORD, RESULT] {
       def create(d : ModelDescriptor) : Option[Model[RECORD, RESULT]]
       def restore(bytes : Array[Byte]) : Model[RECORD, RESULT]
Lightbend
```

Corresponding Scala code that could be generated from the description, although we hand code this logic in the examples and exercises.

Production Concern: Monitoring

Model monitoring should provide information about usage, behavior, performance and lifecycle of the deployed models

```
case class ModelToServeStats(
                                          // Scala example
                                         // Model name
name: String,
description: String,
                                         // Model descriptor
modelType: ModelDescriptor.ModelType, // Model type
                                         // Start time of model usage
since: Long,
                                         // Number of records scored
usage: Long = 0,
duration: Double = 0.0,
                                         // Time spent on scoring
min: Long = Long.MaxValue,
                                         // Min scoring time
max: Long = Long.MinValue
                                         // Max scoring time
                                         We'll return to production concerns...
  Lightbend
```



Note the "ad hoc" part. It's for times when a "normal" stream of output data isn't the best fit, e.g., only periodic updates are needed, you really want to support a range of "impromptu" (ad hoc) queries, etc.

Kafka Streams and Flink have built-in support for this and it's being added to Spark Streaming. We'll show how to use other Akka features to provide the same ability in a straightforward way for Akka Streams.

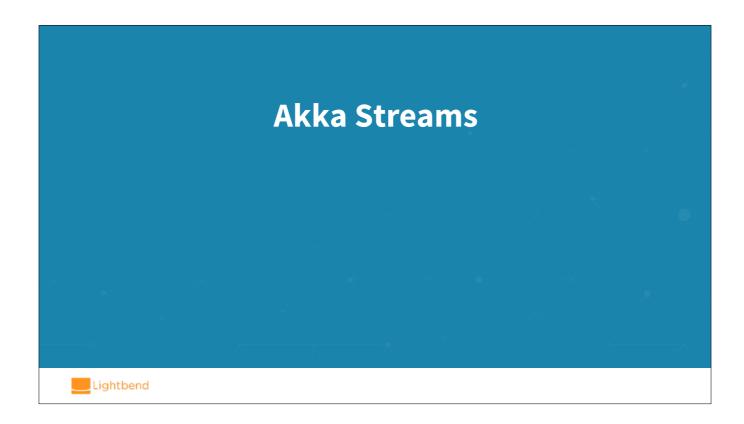
Example used in this tutorial

Throughout the rest of tutorial we use models based on *Wine quality* data for training and scoring.

- The data is publicly available at https://www.kaggle.com/vishalyo990/prediction-of-quality-of-wine/data.
- A great notebook describing this data and providing several machine learning algorithms for this problem: https://www.kaggle.com/vishalyo990/prediction-of-quality-of-wine/notebook

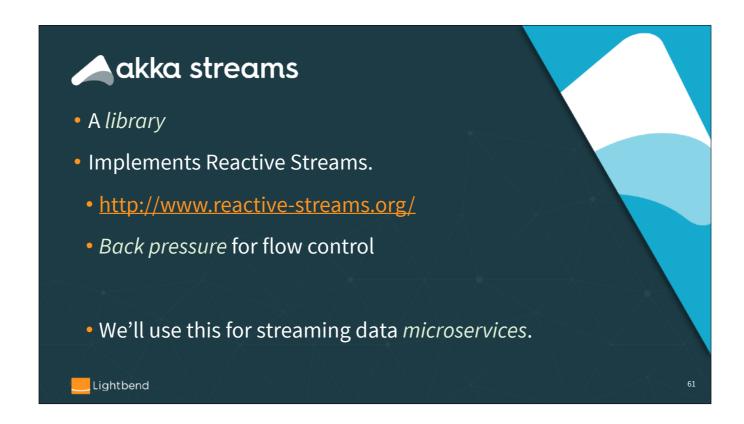
Lightbend

5

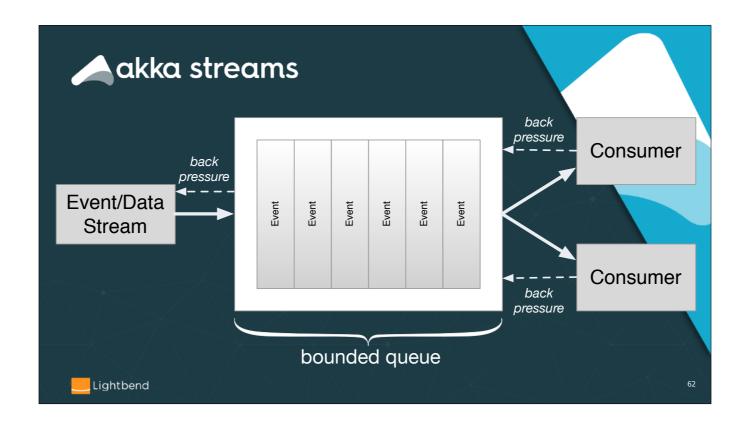




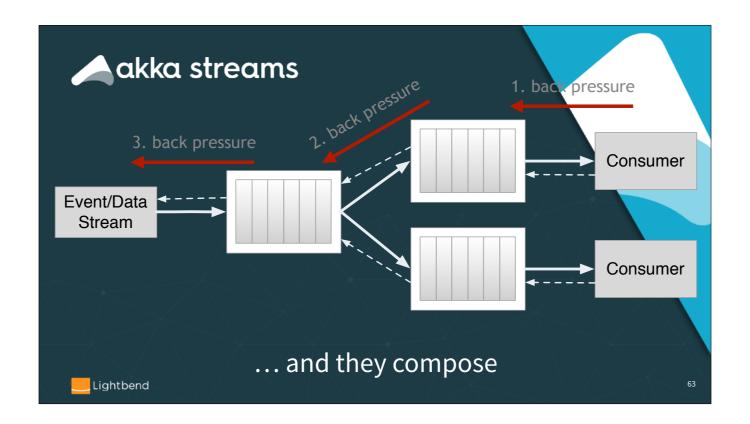
We provide two implementations. You could generalize the second approach (async calls) to invoke an external service. We won't provide examples of this option, but return later with some additional considerations about it.



See this website for details on why back pressure is an important concept for reliable flow control, especially if you don't use something like Kafka as your "near-infinite" buffer between services.



Bounded queues are the only sensible option (even Kafka topic partitions are bounded by disk sizes), but to prevent having to drop input when it's full, consumers signal to producers to limit flow. Most implementations use a push model when flow is fine and switch to a pull model when flow control is needed.



And they compose so you get end-to-end back pressure.



Rich, mature tools for the full spectrum of microservice development.



- A very simple example to get the "gist":
- Calculate the factorials for n = 1 to 10

Lightbend

65

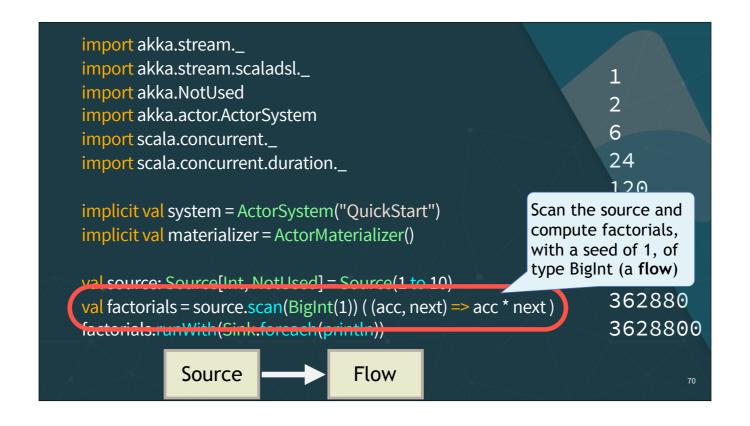
```
import akka.stream._
import akka.stream.scaladsl._
import akka.NotUsed
import akka.actor.ActorSystem
                                                                 6
import scala.concurrent._
                                                                  24
import scala.concurrent.duration._
                                                                  120
implicit val system = ActorSystem("QuickStart")
                                                                 720
implicit val materializer = ActorMaterializer()
                                                                  5040
                                                                 40320
val source: Source[Int, NotUsed] = Source(1 to 10)
                                                                 362880
val factorials = source.scan(BigInt(1)) ( (acc, next) => acc * next )
                                                                 3628800
factorials.runWith(Sink.foreach(println))
```

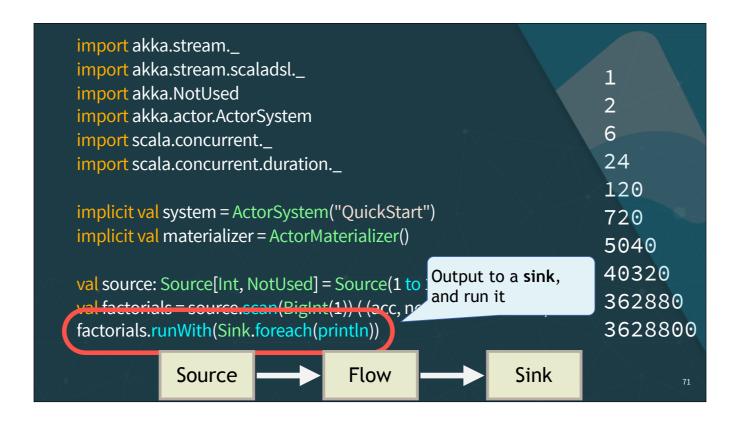
This example is in akkaStreamsModelServer/simple-akka-streams-example.sc

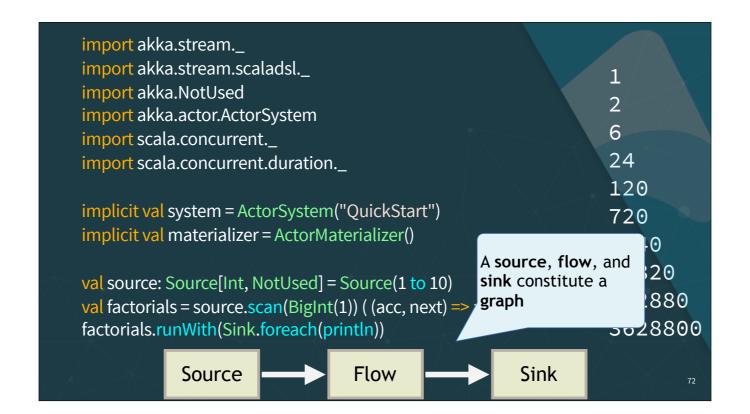
```
import akka.stream._
import akka.stream.scaladsl._
                                      Imports!
import akka.NotUsed
import akka.actor.ActorSystem
                                                                 6
import scala.concurrent._
                                                                 24
import scala.concurrent.duration._
                                                                 120
implicit val system = ActorSystem("QuickStart")
                                                                 720
implicit val materializer = ActorMaterializer()
                                                                 5040
                                                                 40320
val source: Source[Int, NotUsed] = Source(1 to 10)
                                                                 362880
val factorials = source.scan(BigInt(1)) ( (acc, next) => acc * next )
                                                                 3628800
factorials.runWith(Sink.foreach(println))
```

import akka.stream._ import akka.stream.scaladsl._ import akka.NotUsed import akka.actor.ActorSystem import scala.concurrent._ import scala.concurrent.duration._ Initialize and specify now the stream is implicit val system = ActorSystem("QuickStart") "materialized" implicit val materializer = ActorMaterializer() 5040 40320 val source: Source[Int, NotUsed] = Source(1 to 10) 362880 val factorials = source.scan(BigInt(1)) ((acc, next) => acc * next) 3628800 factorials.runWith(Sink.foreach(println))

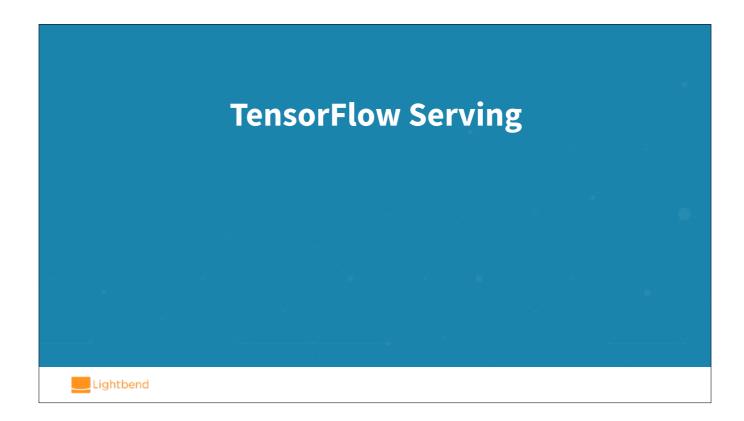
import akka.stream._ import akka.stream.scaladsl._ import akka.NotUsed import akka.actor.ActorSystem import scala.concurrent._ import scala.concurrent.duration._ Create a source of Ints. Second type implicit val system = ActorSystem("QuickStart") represents a hook used for "materialization" implicit val materializer = ActorMaterializer() not used here 40320 val source: Source[Int, NotUsed] = Source(1 to 10) 362880 val factorials - source.scan(Bigint(1)) ((acc, next) -> acc * next) factorials.runWith(Sink.foreach(println)) 3628800 Source

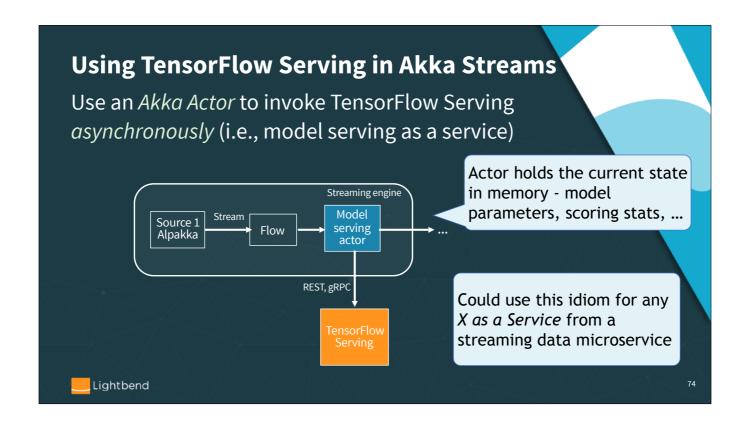






The core concepts are sources and sinks, connected by flows. There is the notion of a Graph for more complex dataflows, but we won't discuss them further





To summarize, we are using Akka Streams to write our streaming data pipeline and calling a custom Akka Actor to handle invoking of the scoring service.

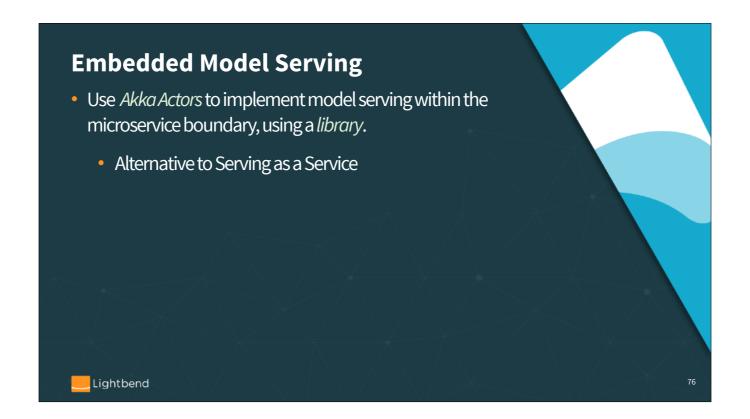
Use the same routing layer idiom: an actor that will implement model serving for a specific model (based on some key) and route messages appropriately to the external service. This way our system can serve models in parallel.

TensorFlow Serving

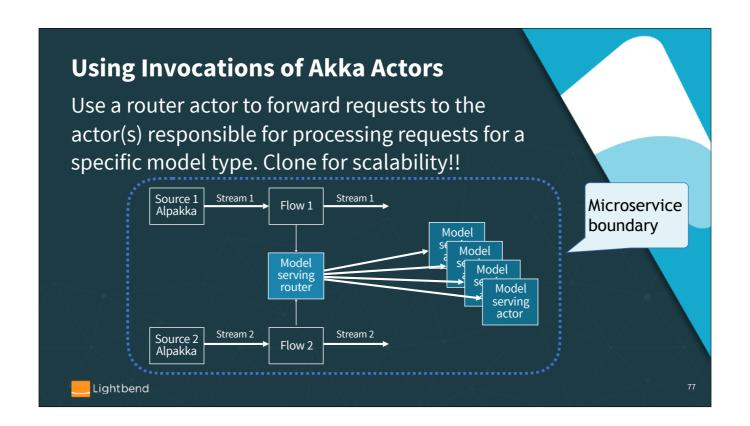
Code time

- Let's look at TF Serving first. We'll use it with microservices shortly.
- Open the example code project
- We'll walk through the project at a high level
- Familiarize yourself with the *tensorflowserver* code
- Load and start the TensorFlow model serving Docker image
- See <u>Using TensorFlow Serving</u> in the README
- Try the implementation and see if you have any questions

Lightbend



We provide both kinds of implementations. You could generalize the second approach (async calls) to invoke an external service. We won't provide examples of this option, but return later with some additional considerations about it.



Create a routing layer: an actor that will implement model serving for a specific model (based on some key) and route messages appropriately. This way our system can serve models in parallel.



Custom stage is an elegant implementation but not scale well to a large number of models. Although a stage can contain a hash map of models, all of the execution will be happening at the same place

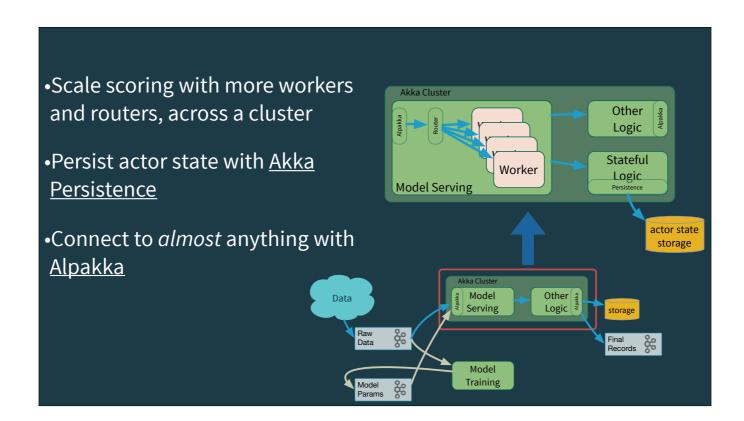


- Implements *queryable state*
- •Curl or open in a browser:

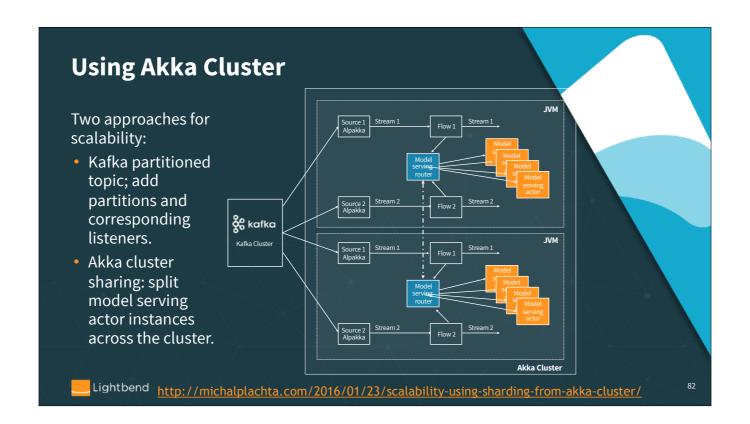
http://localhost:5500/models
http://localhost:5500/state/wine







Here's our streaming microservice example adapted for Akka Streams. We'll still use Kafka topics in some places and assume we're using the same implementation for the "Model Training" microservice. Alpakka provides the interface to Kafka, DBs, file systems, etc. We're showing two microservices as before, but this time running in Akka Cluster, with direct messaging between them. We'll explore this a bit more after looking at the example code.



A great article http://michalplachta.com/2016/01/23/scalability-using-sharding-from-akka-cluster/ goes into a lot of details on both implementation and testing



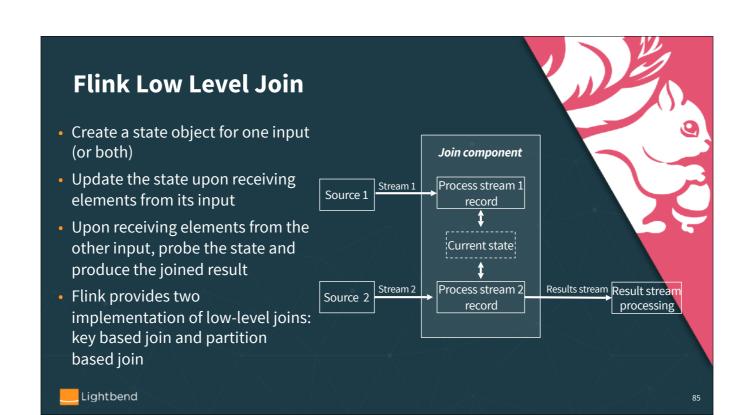
Same sample use case, now with Kafka Streams

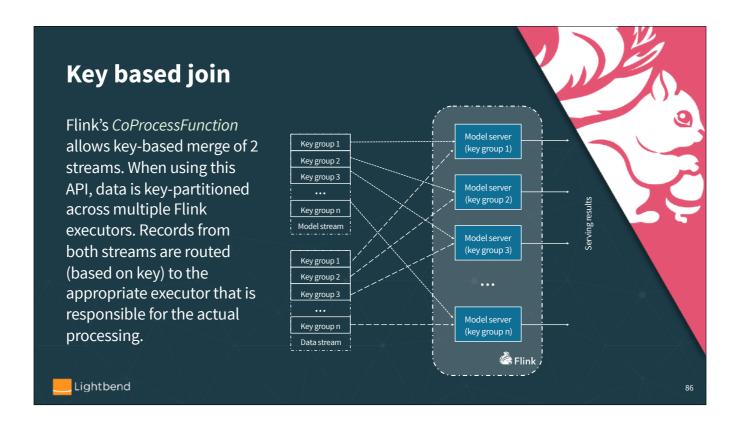


<u>Apache Flink</u> is an open source stream-processing engine (SPE) that provides the following:

- Scales to thousands of nodes.
- Provides checkpointing and save-pointing facilities that enable fault tolerance and the ability to restart without loss of accumulated state.
- Provides queryable state support, which minimizes the need for external databases for external access to the state.
- Provides window semantics, enabling calculation of accurate results, even in the case of out-of-order or late-arriving data.

Lightbend

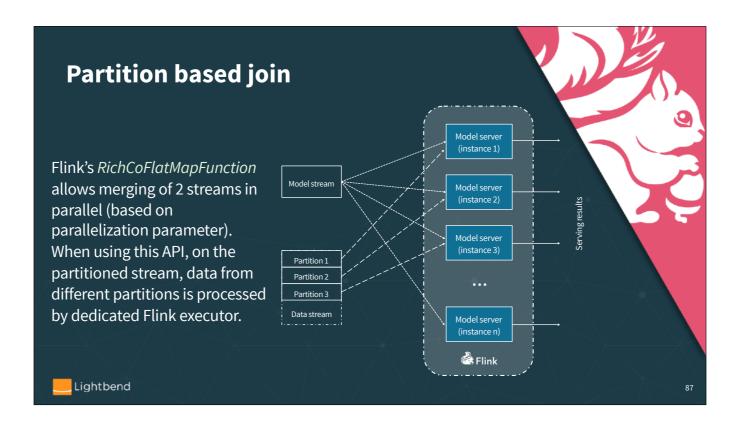




The main characteristics of this approach:

- · Distribution of execution is based on key
- Individual models' scoring is implemented by a separate executor (a single executor can score multiple models), which means that scaling Flink leads to a better distribution of individual models and consequently better parallelization of scorings.
- A given model is always scored by a given executor, which means that depending on the data type distribution of input records, this approach can lead to "hot" executors

Based on this, key-based joins are an appropriate approach for the situations when it is necessary to score multiple data types with relatively even distribution.



Here are the main characteristics of this approach:

• The same model can be scored in one of several executors based on the partitioning of the data streams, which means that scaling of Flink (and input data partitioning) leads to better scoring throughput.

• Because the model stream is broadcast to all model server instances, which operate independently, some race conditions in the model update can exist, meaning that at the point of the model switch, some model jitter (models can be updated at different instances, so for some short period of time different input records can be served by different models) can occur.

Based on these considerations, using global joins is an appropriate approach for the situations when it is necessary to score with one or a few models under heavy data load.

Flink Example

Code time

- 1. Run the *client* project (if not already running)
- 2. Explore and run *flinkServer* project
 - a. ModelServingKeyedJob implements keyed join
 - b. ModelServingFlatJob implements partitioned join

Lightbend

Apache Spark Structured Streaming

Lightbend

Same sample use case, now with Kafka Streams

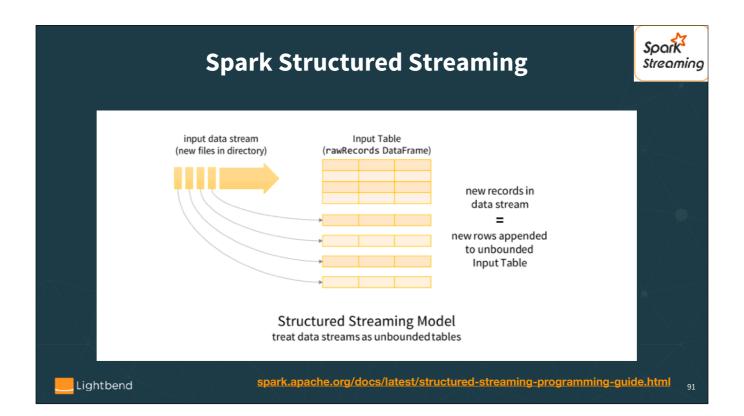
Spark Structured Streaming



<u>Apache Spark</u> Structured Streaming is a scalable and fault-tolerant stream processing engine built on the Spark SQL engine.

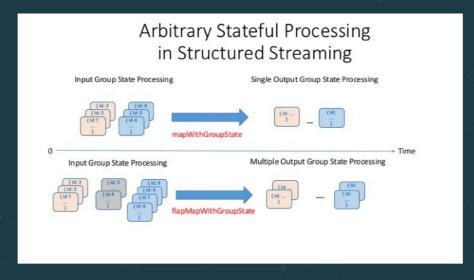
- •Scales to thousands of nodes.
- •Express your streaming computation the same way you would express a batch SQL computation on static data:
- •The Spark SQL engine will take care of running it incrementally and continuously. It updates results as streaming data continues to arrive.

Lightbend



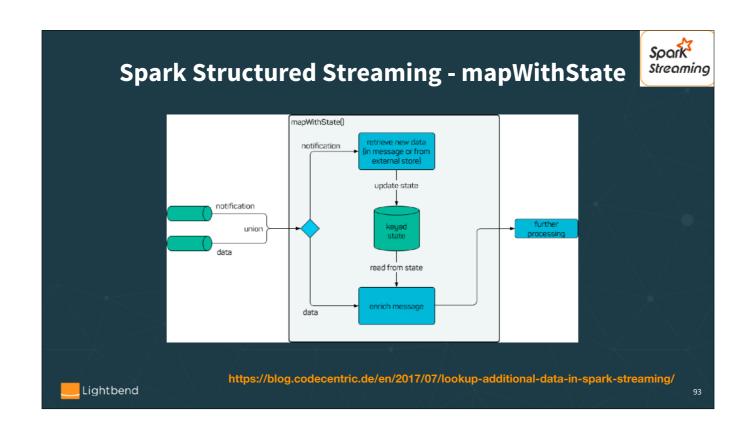
Spark Structured Streaming - State

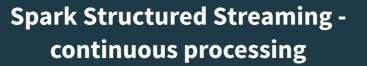




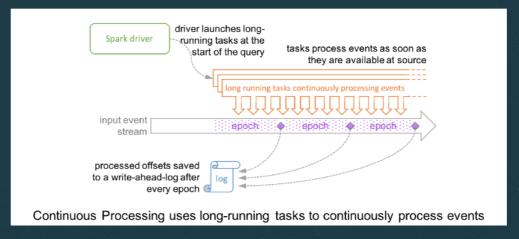
https://databricks.com/blog/2017/10/17/arbitrary-stateful-processing-in-apache-sparks-structured-streaming.html

Lightbend

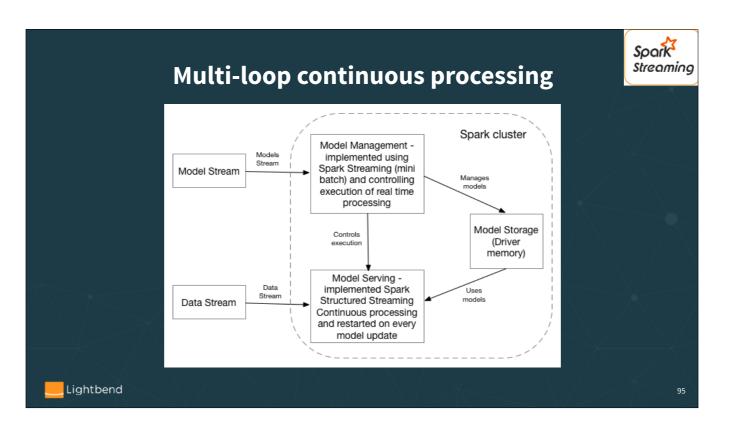








https://databricks.com/blog/2018/03/20/low-latency-continuous-processing-mode-in-structured-streaming-in-apache-spark-2-3-0.html



Spark Example



Code time

- 1. Run the *client* project (if not already running)
- 2. Explore and run the *sparkServer* project
 - a. SparkStructuredModelServer uses mapWithState.
 - b. SparkStructuredStateModelServer implements multiloop approach

Lightbend

Comparing Implementations

- 1. Akka Streams with Akka is a library providing great flexibility for implementations and deployments, but requires custom implementations for scaling and failover.
- More flexibility, but more responsibility (i.e., do it yourself)
- 2. Flink and Spark Streaming are stream-processing engines (SPE) that automatically leverage cluster resources for scaling and failover. Computations are a set of operators and they handle execution parallelism for you, using different threads or different machines.
- Less flexibility, but less responsibility

Lightbend

97

We know Akka Streams best, but other streaming libraries, like Kafka Streams, could be used, too. Be sure they have the flexibility you need for the model serving patterns we discussed! If not, they are JUST libraries, so you can use other libraries to provide missing functionality. We've even used Akka with Kafka Streams...

The sub-bullets summarize the tradeoffs.

Spark vs Flink

- 1. Flink: iterations are executed as cyclic data flows; a program (with all its operators) is scheduled just once and the data is fed back from the tail of an iteration to its head. This allows Flink to keep all additional data locally.
- 2. *Spark*: each iteration is a new set of tasks/operators scheduled and executed. Each iteration operates on the result of the previous iteration which is held in memory. For each new execution cycle, the results have to be moved to the new execution processes.

Lightbend

https://hal.inria.fr/hal-01347638v2/document

Spark vs Flink

- 1. In *Flink*, all additional data is kept locally, so arbitrarily complex structures can be used for its storage (although serializers are required for checkpointing). The serializers are only invoked out of band.
- 2. In *Spark*, all additional data is stored external to each minibatch, so it has to be marshaled/unmarshaled for every minibatch (for *every message* in continuous execution) to make this data available.
- 3. Spark Structured Streaming is based on SQL data types, which makes data storage even more complex (JVM vs. SQL datatyping).

Lightbend

• Hidden technical debt in machine learning systems

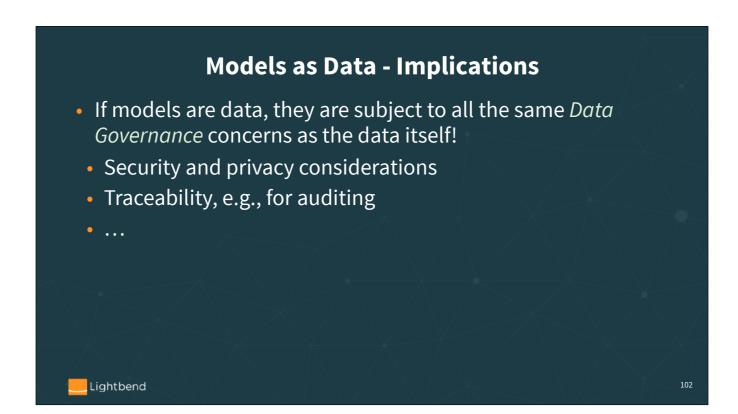
- Model serving patterns
 - Embedding models as code
 - Models as data
 - External services
 - Dynamically controlled streams
- Additional production concerns for model serving
- Wrap up

Lightbend

Additional Production Concerns for Model Serving

- Production implications for *models* as data
- Software process concerns, e.g., CI/CD
- Another pattern: speculative execution of models

Lightbend



We'll just discuss these two aspects of the larger topic of data governance

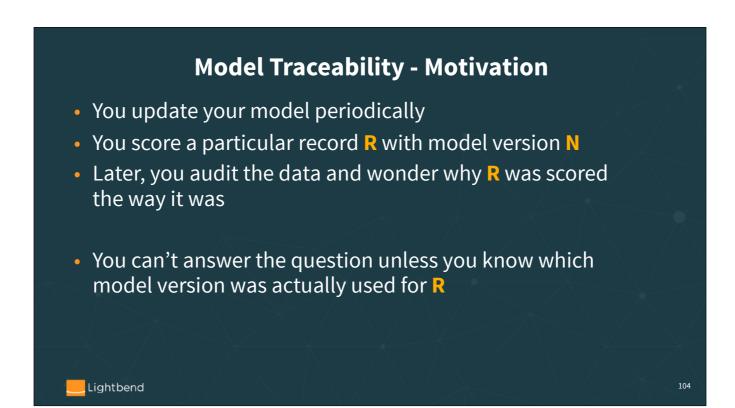
Security and Privacy Considerations

- Models are intellectual property
- So controlled access is required
- How do we preserve privacy in model-training, scoring, and other data usage?
- See these papers and articles on privacy preservation

Lightbend

10

The "papers" link is to a Google Scholar search with several papers on privacy preserving techniques, like differential privacy.



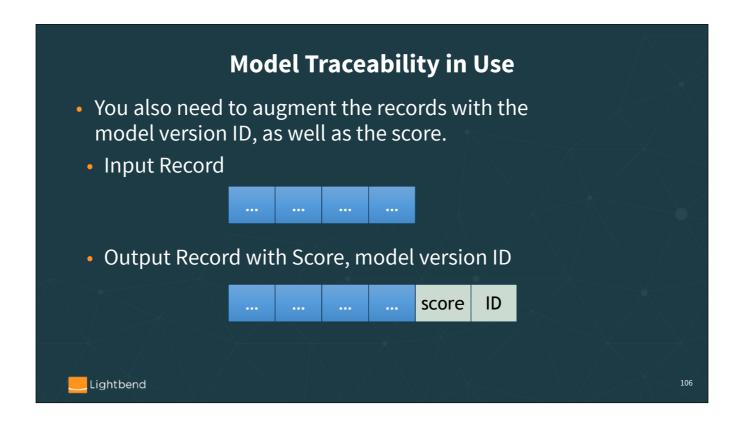
We mentioned this example before. "Explainability" is an important problem in Deep Learning; knowing why the model produced the results it produced.

Model Traceability Requirements

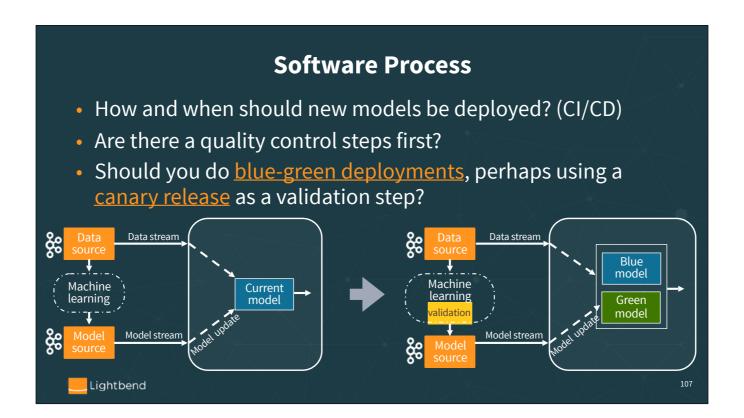
- A model repository
- Information stored for each model instance, e.g.:
- Version ID
- Name, description, etc.
- Creation, deployment, and retirement dates
- Model parameters
- Quality metrics

•

Lightbend



You might be tempted to avoid this extra data; don't you know the "deployment" date? This usually isn't accurate enough to know which records were on the boundary of switchover, due to timestamp granularity, clock skew, etc.



Just the tip of the iceberg...

The blue-green deployments leads to our final topic...

Speculative Execution

According to Wikipedia, speculative execution is an **optimization** technique, where:

- The system performs work that may not be needed, before it's known if it will be needed.
- If and when it *is* needed, we don't have to wait.
- The results are discarded if they aren't needed.

Lightbend

Speculative Execution

- Provides more concurrency if extra resources are available.
- Used for:
- branch prediction in pipelined processors,
- value prediction for exploiting value locality,
- prefetching instructions and files,
- etc.

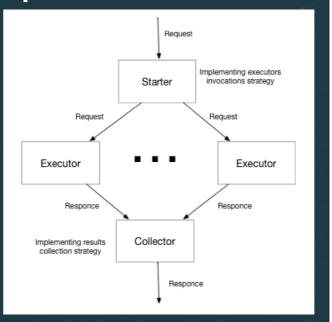
Why not use it with machine learning??

Lightbend

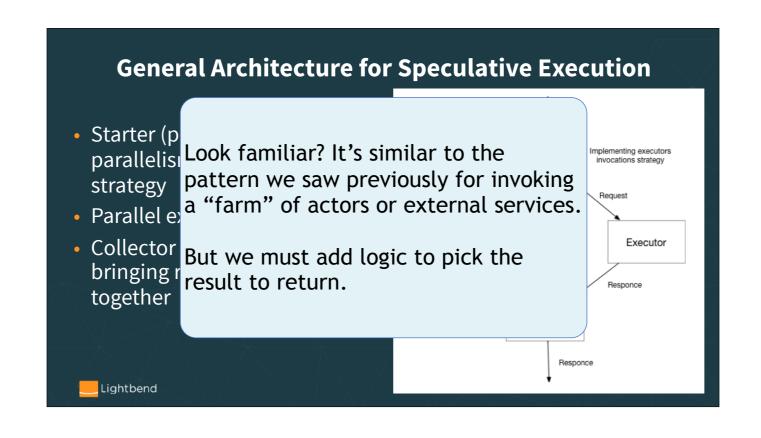
109

General Architecture for Speculative Execution

- Starter (proxy) controls parallelism and invocation strategy
- Parallel execution by executors
- Collector responsible for bringing results from executors together







Use Case - Guaranteed Execution Time

- I.e., meet a tight latency SLA
- Run several models:
- A smart model, but takes time *T1* for a given record
- A "less smart", but faster model with a fixed upper-limit on execution time, with *T2* << *T1*
- If timeout (latency budget) *T* occurs, where *T2 < T < T1*, return the less smart result
- But if *T1* < *T*, return that smarter result
- (Is it clear why T2 < T < T1 is required?)

Lightbend

11

Because the timeout T has to be long enough that the fast model has time to finish, so T must be longer than T2, and this is only useful if T1 is often longer than T, our latency window.

This technique is "speculative", because we try both models, taking a compromise result if the good result takes too long to return.

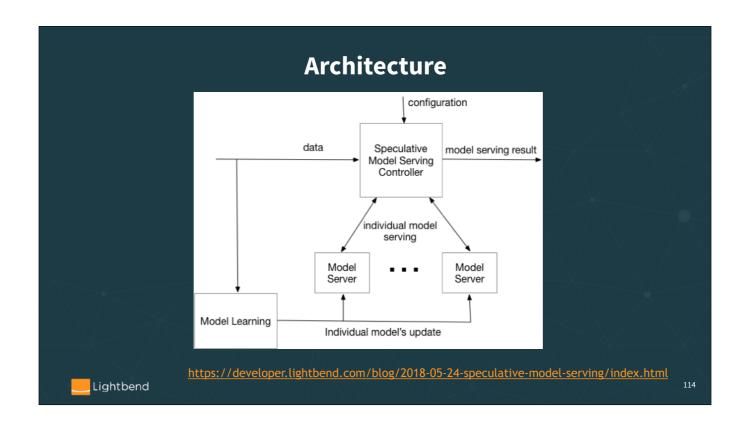
Use Case - Ensembles of Models

- Consensus-based model serving
- *N* models (*N* is *odd*)
- Score with all of them and return the *majority* result
- Quality-based model serving
- N models using the same *quality metric*
- Pick the result for a given record with the best quality score
- Similarly for more sophisticated <u>boosting</u> and <u>bagging</u> systems

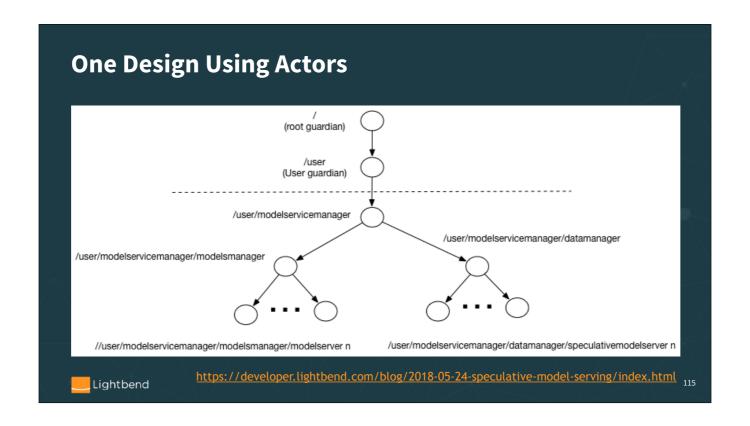


113

Consensus here is a majority vote system, a fairly crude *ensemble method*, while boosting and bagging are more sophisticated ensemble systems. Quality is roughly speaking an ensemble system, but here the models are treated independently, all with a quality metric, and the best score, based on that metric, is chosen.



This blog post provides more information on this technique.



The path-like strings are the Akka way of defining a hierarchy of actors and provide an abstract way to reference an actor that doesn't require you to know the actual process and machine where it's running.

• Hidden technical debt in machine learning systems

- Model serving patterns
 - Embedding models as code
 - Models as data
 - External services
 - Dynamically controlled streams
- Additional production concerns for model serving
- Wrap up

Lightbend

116

Recap

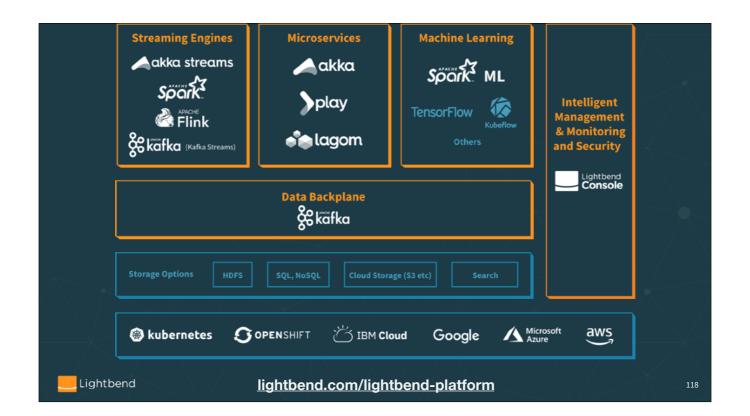
- Model serving is one small(-ish) part of the whole ML pipeline
- Use logs (e.g., Kafka) to connect most services
- Models as data provides the most flexibility
- Model serving can be implemented in "general" microservices, like *Akka Streams*, or data systems, like *Flink* and *Spark*
- Model serving can be *in-process* (embedded library) or an *external service* (e.g., TensorFlow Serving)
- Production concerns include integration with your CI/CD pipeline and data governance

Lightbend

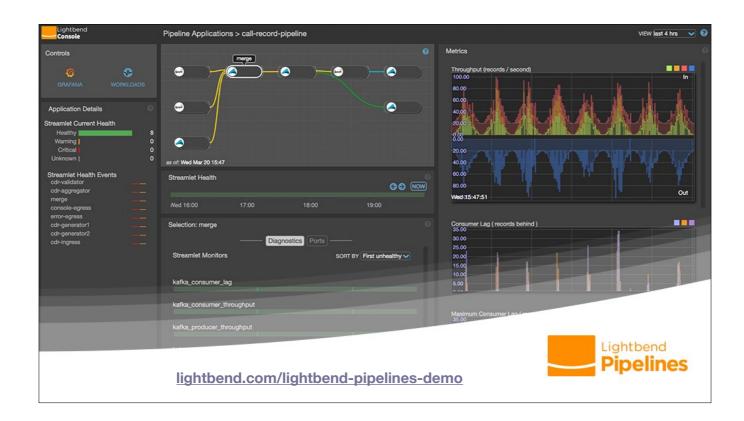
117

Some of the main points we discussed.

See this blog post for more about when to use Akka Streams and when not to use it. https://medium.com/@unmeshvjoshi/choosing-a-stream-processing-framework-spark-streaming-or-kafka-streams-or-alpakka-kafka-fa0422229b25



<marketing>Lightbend has been thinking about the development of microservices and streaming data applications for a while. Lightbend Platform is a commercial distribution of open-source and commercial tools to accelerate development, deployment, and management of these applications.

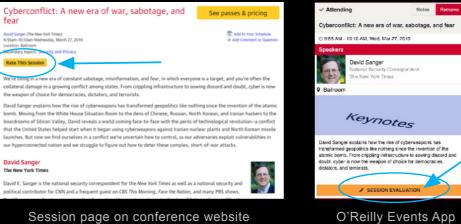


<marketing>In particular, Lightbend Pipelines is a new capability in Lightbend Platform for writing, deploying, monitoring, and managing streaming data pipelines.
marketing>

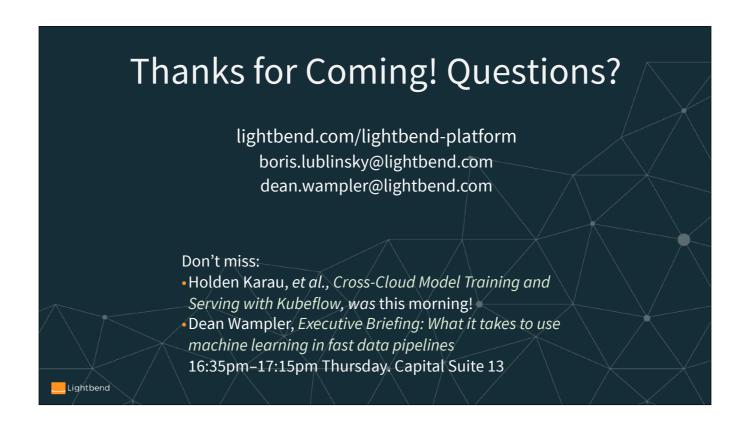


We can only touch on architectures, design principles, and implementation idioms in this tutorial. For more extensive coverage, see these free reports written by Boris and Dean, which are published by O'Reilly.

Rate today's session



Session page on conference website



Thank you! Please check out Dean's other session. Also, Holden, et al. had a tutorial this morning specifically on Kubeflow. You'll be able to watch the video on O'Reilly's Safari. Finally, please check out our Lightbend Platform for commercial options for building and running microservices (with or without ML) using Kafka, Spark, Flink, Akka Streams, and Kafka Streams.