

# Balancing the Dataset - AWS Glue Job Script (& III)

## Write samples in S3

```
# Build a Spark DataFrame that is the union of all the samples
balanced_df = samples[0]

for sample in samples[1:]:
    balanced_df = balanced_df.union(sample)

# Convert the DataFrame to a Glue DynamicFrame
balanced = DynamicFrame.fromDF(balanced_df, glueContext, "balanced")

# Write the data on the target bucket using parquet format
sampled_data_sink = glueContext.write_dynamic_frame.from_options(
    frame = balanced,
    connection_type='s3',
    connection_options={"path": target, "partitionKeys": ["product_category"]},
    format="parquet")
```

# SageMaker BlazingText

Two modes:

- Unsupervised
  - Highly optimized implementation of the **Word2vec**
  - Used for converting words to vectors (aka ***word embeddings***)
- Supervised
  - Extends the **fastText** text classifier
  - Used for multi class/label text classification

# Target Format - BlazingText

- Single preprocessed text file
- Space separated tokens
- Single sentence per line
- Labels alongside the sentence
- A label is a word prefixed by the string `_label_`
- Can use training and validation channels

## TODOs:

- Select only the used fields
- Tokenize review body and convert it to a space separated string and prepend the string with the product category label
- Split the dataset into training, validation, and test subsets
- Write each subset to a single object in S3

# Preparing the Dataset - AWS Glue Job Script (I)

Select only the used fields

```
# Read from the data source
datasource = glueContext.create_dynamic_frame.from_catalog(
    database = database,
    table_name = table)

# Select only the fields that are going to be used
select = SelectFields.apply(
    frame = datasource,
    paths = ["product_category", "review_body"])
```

# Preparing the Dataset - AWS Glue Job Script (II)

Tokenize review body and convert it to a space separated string

```
# Transform the reviews text by applying the tokenize function to each row in the DynamicFrame
tokenized = Map.apply(frame = select, f = tokenize, transformation_ctx = "tokenized")

def tokenize(dynamicRecord):
    category = dynamicRecord['product_category'].lower()
    dynamicRecord['product_category'] = '__label__'+dynamicRecord['product_category'].lower()
    dynamicRecord['review'] = transform_review_body(dynamicRecord.get('review_body', ''))
    return dynamicRecord

def transform_review_body(review_body):
    from nltk.tokenize import TweetTokenizer
    tknzs = TweetTokenizer()
    body = tknzs.tokenize(remove_tags(review_body.lower()))

    return ' '.join(body)
```

# Preparing the Dataset - AWS Glue Job Script (III)

Split the dataset into training, validation, and test subsets

```
# Split the sample into train, test, and validation sets. A Spark DataFrame is needed for th  
df = tokenized.toDF()  
train, validation, test = df.randomSplit(weights = [.6, .2, .2], seed = 42)
```

# Preparing the Dataset - AWS Glue Job Script (IV)

Write each subset to a single object in S3

```
# Repartition the data frames to store each set into a single file and convert to DynamicFrame
train_set = DynamicFrame.fromDF(train.repartition(1), glueContext, "train")
validation_set = DynamicFrame.fromDF(validation.repartition(1), glueContext, "validation")
test_set = DynamicFrame.fromDF(test.repartition(1), glueContext, "test")

# Write each set with a data sink
train_datasink = glueContext.write_dynamic_frame.from_options(
    frame = train_set,
    connection_type = "s3",
    connection_options = {"path": "{}/train".format(target)},
    format = "csv",
    format_options = {"separator": " ", "writeHeader": False, "quoteChar": "-1" })
```

# Using the generic Amazon SageMaker Estimator

- Part of the SageMaker Python SDK
- Needs to be configured with:
  - Container for BlazingText algorithm
  - Training and validation channels
  - IAM Role
  - Training instance configuration
  - Hyperparameters
  - Output location for the training artifacts

# jupyter reviews-classification-hpt Last Checkpoint: a day ago (autosaved)



File Edit View Insert Cell Kernel Widgets Help

Trusted

conda\_python3



from sagemaker python SDK.

```
In [15]: import sagemaker
import json
import boto3

sess = sagemaker.Session()

role = sagemaker.get_execution_role()
print(role) # This is the role that SageMaker would use to leverage AWS resources (S3, CloudWatch, etc)

bucket = 'sagemaker-reviews-demo'

job_run_id = 'jr_17371f644effb76dfb62c4a4bd3f305e09ee02f87a7723a5c83d83a5e2e7579b'

prefix = 'transformed/' + job_run_id #Replace with the prefix under which you want to store the transformed data

arn:aws:iam::XXXXXXXXXXXX:role/service-role/AWSGlueServiceSageMakerNotebookRole-TextClassification
```

You are going to need the nltk library, which is not included in the notebook instance, so let's use pip to install it (this is only needed once).

```
In [2]: !pip install nltk
```

# jupyter reviews-classification-hpt Last Checkpoint: a day ago (autosaved)



File Edit View Insert Cell Kernel Widgets Help

Trusted

conda\_python3



```
prefix = 'transformed/' + job_run_id #Replace with the prefix under which you want to store the  
arn:aws:iam:::role/service-role/AWSGlueServiceSageMakerNotebookRole-TextClassification
```

You are going to need the nltk library, which is not included in the notebook instance, so let's use pip to install it (this is only needed once).

In [2]: `!pip install nltk`

```
Requirement already satisfied: nltk in /home/ec2-user/anaconda3/envs/python3/lib/python3.6/site-packages (3.3)
Requirement already satisfied: six in /home/ec2-user/anaconda3/envs/python3/lib/python3.6/site-packages (from nltk) (1.11.0)
You are using pip version 10.0.1, however version 18.1 is available.
You should consider upgrading via the 'pip install --upgrade pip' command.
```

In [3]: `import nltk  
nltk.download('punkt')`

```
[nltk_data] Downloading package punkt to /home/ec2-user/nltk_data...
[nltk_data]  Unzipping tokenizers/punkt.zip.
```

Out[3]: True

# jupyter reviews-classification-hpt Last Checkpoint: a day ago (autosaved)



File Edit View Insert Cell Kernel Widgets Help

Trusted

conda\_python3



Out[3]: True

## Data channels

Amazon SageMaker requires a channel for reading training data and another for validation data. To set this up, you need to create the `sagemaker.session.s3_input` objects from the data channels. These objects are then put in a simple dictionary, which the algorithm consumes.

```
In [4]: train_channel = prefix + '/train/'
validation_channel = prefix + '/validation/'

s3_train_data = 's3://{}{}'.format(bucket, train_channel)
s3_validation_data = 's3://{}{}'.format(bucket, validation_channel)

train_data = sagemaker.session.s3_input(s3_train_data, distribution='FullyReplicated',
                                       content_type='text/plain', s3_data_type='S3Prefix')
validation_data = sagemaker.session.s3_input(s3_validation_data, distribution='FullyReplicated',
                                             content_type='text/plain', s3_data_type='S3Prefix')

data_channels = {'train': train_data, 'validation': validation_data}
```

Next you need to setup an output location at S3, where the model artifact will be dumped. These artifacts are also the output of the algorithm's training job.

## jupyter reviews-classification-hpt Last Checkpoint: a day ago (autosaved)



File Edit View Insert Cell Kernel Widgets Help

Trusted

conda\_python3



the algorithm's training job.

```
In [5]: s3_output_location = 's3://{}{}/output/{}'.format(bucket, job_run_id)
```

## Training

Now that you are done with all the setup that is needed, you are ready to train our object detector. To begin, create a `sageMaker.estimator.Estimator` object. This estimator will launch the training job.

```
In [6]: region_name = boto3.Session().region_name
```

```
In [7]: container = sagemaker.amazon.amazon_estimator.get_image_uri(region_name, "blazingtext", "latest")
print('Using SageMaker BlazingText container: {} {}'.format(container, region_name))

Using SageMaker BlazingText container: 685385470294.dkr.ecr.eu-west-1.amazonaws.com/blazingtext:latest (eu-west-1)
```

## Training the BlazingText model for supervised text classification

Similar to the original implementation of [Word2Vec](#), SageMaker BlazingText provides an efficient implementation of the continuous bag-of-words (CBOW) and skip-gram architectures using Negative Sampling, on CPUs and additionally on GPU[s].

 jupyter reviews-classification-hpt Last Checkpoint: a day ago (autosaved)

[File](#) [Edit](#) [View](#) [Insert](#) [Cell](#) [Kernel](#) [Widgets](#) [Help](#)

Trusted

conda\_python3



Modes	cbow (supports subwords training)	skipgram (supports subwords training)	batch_skipgram	supervised
Single CPU instance	✓	✓	✓	✓
Single GPU instance	✓	✓	✓	✓ (Instance with 1 GPU only)
Multiple CPU instances			✓	

Now, let's define the SageMaker Estimator with resource configurations and hyperparameters to train Text Classification on the prepared dataset, using "supervised" mode on a p3.2xlarge instance.

```
In [8]: bt_model = sagemaker.estimator.Estimator(container,
                                                role,
                                                train_instance_count=1,
                                                train_instance_type='ml.c5.4xlarge',
                                                train_volume_size = 30,
                                                train_max_run = 360000,
                                                input_mode= 'File',
                                                output_path=s3_output_location,
                                                sagemaker_session=sess)
```

Please refer to [algorithm documentation]

([https://docs.aws.amazon.com/sagemaker/latest/dg/blazingtext\\_hyperparameters.html](https://docs.aws.amazon.com/sagemaker/latest/dg/blazingtext_hyperparameters.html)) for the complete list of hyperparameters.



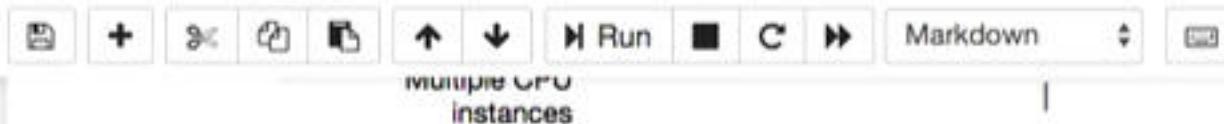
jupyter reviews-classification-hpt Last Checkpoint: a day ago (autosaved)



File Edit View Insert Cell Kernel Widgets Help

Trusted

| conda\_python3 0



Now, let's define the SageMaker Estimator with resource configurations and hyperparameters to train Text Classification on the prepared dataset, using "supervised" mode on a p3.2xlarge instance.

Please refer to [algorithm documentation].

([https://docs.aws.amazon.com/sagemaker/latest/dg/blazingtext\\_hyperparameters.html](https://docs.aws.amazon.com/sagemaker/latest/dg/blazingtext_hyperparameters.html)) for the complete list of hyperparameters.

# jupyter reviews-classification-hpt Last Checkpoint: a day ago (autosaved)



File Edit View Insert Cell Kernel Widgets Help

Trusted

conda\_python3



Run

Cell

Markdown



```
train_instance_type='ml.c5.4xlarge',
train_volume_size = 30,
train_max_run = 360000,
input_mode= 'File',
output_path=s3_output_location,
sagemaker_session=sess)
```

Please refer to [algorithm documentation]

([https://docs.aws.amazon.com/sagemaker/latest/dg/blazingtext\\_hyperparameters.html](https://docs.aws.amazon.com/sagemaker/latest/dg/blazingtext_hyperparameters.html)) for the complete list of hyperparameters.

In [9]: bt\_model.set\_hyperparameters(mode='supervised',

```
min_epochs=5,
epochs=100,
min_count=2,
learning_rate=0.010817706468703524,
vector_dim=134,
early_stopping=True,
patience=30,
word_ngrams=5)
```

We will use SageMaker hyperparameter tuning to automate the searching process effectively. Specifically, we specify a range, or a list of possible values in the case of categorical hyperparameters, for each of the hyperparameter that we plan to tune. SageMaker hyperparameter tuning will automatically launch multiple training jobs with different hyperparameter settings, and the results of these training jobs are used to find the best hyperparameter settings for our model.



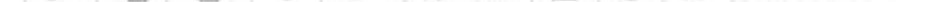
jupyter reviews-classification-hpt Last Checkpoint: a day ago (autosaved)



File Edit View Insert Cell Kernel Widgets Help

Trusted

```
| conda python3 0
```



```
    patience=30,  
    word_ngrams=5)
```

We will use SageMaker hyperparameter tuning to automate the searching process effectively. Specifically, we specify a range, or a list of possible values in the case of categorical hyperparameters, for each of the hyperparameter that we plan to tune. SageMaker hyperparameter tuning will automatically launch multiple training jobs with different hyperparameter settings, evaluate results of those training jobs based on a predefined "objective metric", and select the hyperparameter settings for future attempts based on previous results. For each hyperparameter tuning job, we will give it a budget (max number of training jobs) and it will complete once that many training jobs have been executed.

First we define the ranges for some of the BlaingText supervised mode algorithm

Then define the objective metric:

```
In [11]: objective.metric_name = 'validation:accuracy'
```

# jupyter reviews-classification-hpt Last Checkpoint: a day ago (autosaved)



File Edit View Insert Cell Kernel Widgets Help

Trusted

conda\_python3



We will use SageMaker hyperparameter tuning to automate the searching process effectively. Specifically, we specify a range, or a list of possible values in the case of categorical hyperparameters, for each of the hyperparameter that we plan to tune. SageMaker hyperparameter tuning will automatically launch multiple training jobs with different hyperparameter settings, evaluate results of those training jobs based on a predefined "objective metric", and select the hyperparameter settings for future attempts based on previous results. For each hyperparameter tuning job, we will give it a budget (max number of training jobs) and it will complete once that many training jobs have been executed.

First we define the ranges for some of the BlazingText supervised mode algorithm

```
In [10]: from sagemaker.tuner import IntegerParameter, CategoricalParameter, ContinuousParameter, Hyperparameter_ranges = {'min_count': IntegerParameter(2, 4),  
                           'learning_rate': ContinuousParameter(0.01, 0.08),  
                           'vector_dim': IntegerParameter(100, 200),  
                           'word_ngrams': IntegerParameter(2, 5)}
```

Then define the objective metric

```
In [11]: objective_metric_name = 'validation:accuracy'
```

Now, you'll create a `HyperparameterTuner` object, to which you pass:

- The `BlazingText` estimator created above



https://eu-west-1.console.aws.amazon.com/sagemaker/home?region=eu-west-1#hyper-tuning-jobs/reviews-17371f644eff-181023-0921?region=eu-west-1&region=eu-west-1&tab=bestTrainingJob

aws Services Resource Groups AWS Glue Athena QuickSight S3 Amazon Ireland Support

**Amazon SageMaker**

Dashboard

Notebook

Notebook instances Lifecycle configurations

Training

Training jobs Hyperparameter tuning jobs

Inference

Name	Status	Objective metric	Value
reviews-17371f644eff-181023-0921-009-541efe76	Completed	validation:accuracy	0.6556000113487244

### Best training job hyperparameters

Name	Type	Value
_tuning_objective_metric	FreeText	validation:accuracy
early_stopping	FreeText	True
epochs	FreeText	100
learning_rate	Continuous	0.01
min_count	Integer	2
min_epochs	FreeText	5
mode	FreeText	supervised
patience	FreeText	30
sagemaker_estimator_class_name	FreeText	"Estimator"
sagemaker_estimator_module	FreeText	"sagemaker.estimator"
vector_dim	Integer	183

Feedback English (US) © 2006–2018 Amazon Web Services, Inc. or its affiliates. All rights reserved. Privacy Policy Terms of Use

# Four phases of ML



Discovery



Integration



Training



Deployment

# jupyter reviews-classification-hpt



File Edit View Insert Cell Kernel Widgets Help

Trusted

conda\_python3



```
In [ ]: boto3.client('sagemaker').describe_hyper_parameter_tuning_job(  
    HyperParameterTuningJobName=tuner.latest_tuning_job.job_name)[ 'HyperParameterTuningJobStatus'
```

## Hosting / Inference

Once the training is done, we can deploy the trained model as an Amazon SageMaker real-time hosted endpoint. This will allow us to make predictions (or inference) from the model. Note that we don't have to host on the same type of instance that we used to train. Because instance endpoints will be up and running for long, it's advisable to choose a cheaper instance for inference. **If you trained the model using the hyperparameter tuning, use the console**

```
# In [ ]: tuner.deploy(initial_instance_count = 1, instance_type = 'ml.m4.xlarge', endpoint_name='textclas  
#text_classifier = bt_model.deploy(initial_instance_count = 1,instance_type = 'ml.m4.xlarge')
```

### Use JSON format for inference

BlazingText supports `application/json` as the content-type for inference. The payload should contain a list of sentences with the key as "`instances`" while being passed to the endpoint.

```
In [ ]: with  
sentences = ["disappointed, stopped working due to water after only two weeks",  
            "10 year old loved it!"]
```

## jupyter test\_model Last Checkpoint: 23/10/2018 (autosaved)



File Edit View Insert Cell Kernel Widgets Help

Not Trusted

| conda\_python3 ○



```
# using the same nltk tokenizer that we used during data preparation for training
tokenized_sentences = [' '.join(nltk.word_tokenize(sent)) for sent in sentences]

payload = {"instances": tokenized_sentences}

response = text_classifier.predict(json.dumps(payload))

predictions = json.loads(response)
print(json.dumps(predictions, indent=2))
```

```
[nltk_data] Downloading package punkt to /home/ec2-user/nltk_data...
[nltk_data] Package punkt is already up-to-date!
```

```
[
  {
    "prob": [
      0.39650389552116394
    ],
    "label": [
      "__label__lawn_and_garden"
    ]
  },
  {
    "prob": [
      0.5274955630302429
    ],
    "label": [
      "__label__toys"
    ]
  }
]
```

# jupyter test\_model Last Checkpoint: 23/10/2018 (autosaved)



File Edit View Insert Cell Kernel Widgets Help

Not Trusted | conda\_python3



```
In [15]: payload = {"instances": tokenized_sentences,
                  "configuration": {"k": 3}}
response = text_classifier.predict(json.dumps(payload))
predictions = json.loads(response)
print(json.dumps(predictions, indent=2))
```

```
[
  {
    "prob": [
      0.39650389552116394,
      0.19785909354686737,
      0.1759667694568634
    ],
    "label": [
      "__label__lawn_and_garden",
      "__label__watches",
      "__label__major_appliances"
    ]
  },
  {
    "prob": [
      0.5274955630302429,
      0.31088608503341675,
      0.05716495215892792
    ],
    "label": [
      "__label__lawn_and_garden",
      "__label__watches",
      "__label__major_appliances"
    ]
  }
]
```

# jupyter test\_model Last Checkpoint: 23/10/2018 (autosaved)



File Edit View Insert Cell Kernel Widgets Help

Not Trusted | conda\_python3



```
from sklearn.metrics import confusion_matrix

def plot_confusion_matrix(cm, classes,
                         normalize=False,
                         title='Confusion matrix',
                         cmap=plt.cm.Blues):
    """
    This function prints and plots the confusion matrix.
    Normalization can be applied by setting `normalize=True`.
    """
    if normalize:
        cm = cm.astype('float') / cm.sum(axis=1)[:, np.newaxis]
        print("Normalized confusion matrix")
    else:
        print('Confusion matrix, without normalization')

    print(cm)

    plt.imshow(cm, interpolation='nearest', cmap=cmap)
    plt.title(title)
    plt.colorbar()
    tick_marks = np.arange(len(classes))
    plt.xticks(tick_marks, classes, rotation=45)
    plt.yticks(tick_marks, classes)

    fmt = '.2f' if normalize else 'd'
    thresh = cm.max() / 2.
```

# AWS Glue + Amazon SageMaker Cost

Development endpoint:

- 2 DPUs x 3 hours x \$0.44 DPU-hour = \$2.64

Sampling Job:

- 30 DPUs x 6 min (minimum 10) x \$0.44 DPU-hour = \$2.2

Preparation Job:

- 20 DPUs x 25 min x \$0.44 DPU-hour = \$3.67

SageMaker Notebook instance:

- 1 ml.t2.medium x 3 hours x \$0.05 = \$0.15

Hyperparameter Tuning Job:

- 1 ml.c5.4xlarge x \$1.075 x 2 h 58 min total training = \$3.19

**Total: \$11.85**

# AWS Glue + Amazon SageMaker Cost

Development endpoint:

- 2 DPUs x 3 hours x \$0.44 DPU-hour = \$2.64

Sampling Job:

- 30 DPUs x 6 min (minimum 10) x \$0.44 DPU-hour = \$2.2

Preparation Job:

- 20 DPUs x 25 min x \$0.44 DPU-hour = \$3.67

SageMaker Notebook instance:

- 1 ml.t2.medium x 3 hours x \$0.05 = \$0.15

Hyperparameter Tuning Job:

- 1 ml.c5.4xlarge x \$1.075 x 2 h 58 min total training = \$3.19

**Total: \$11.85**

## AWS Glue + Amazon SageMaker Cost

Development endpoint:

- 2 DPUs x 3 hours x \$0.44 DPU-hour = \$2.64

Sampling Job:

- 30 DPUs x 6 min (minimum 10) x \$0.44 DPU-hour = \$2.2

Preparation Job:

- 20 DPUs x 25 min x \$0.44 DPU-hour = \$3.67

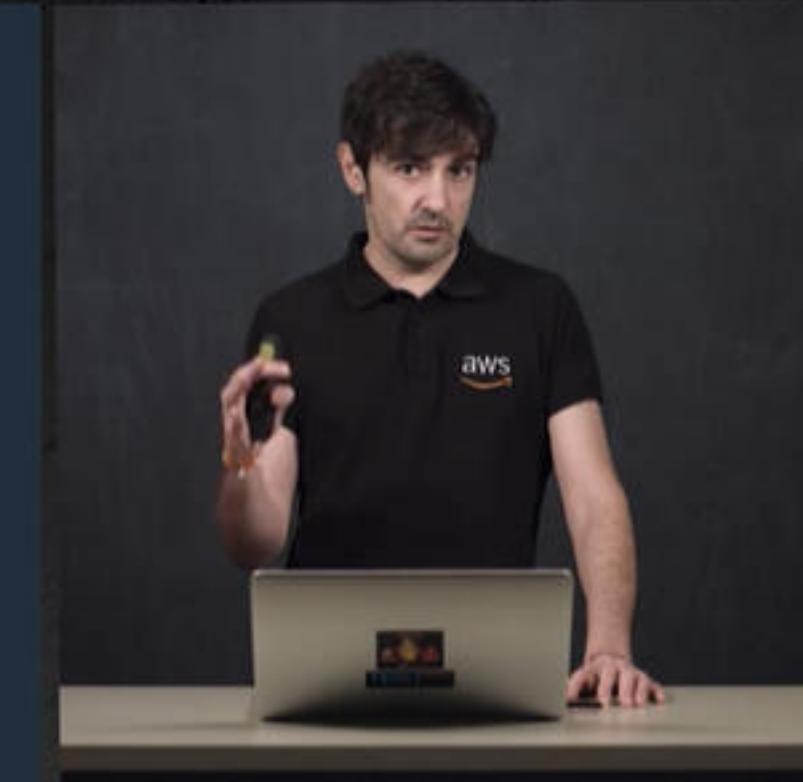
SageMaker Notebook instance:

- 1 ml.t2.medium x 3 hours x \$0.05 = \$0.15

Hyperparameter Tuning Job:

- 1 ml.c5.4xlarge x \$1.075 x 2 h 58 min total training = \$3.19

**Total: \$11.85**



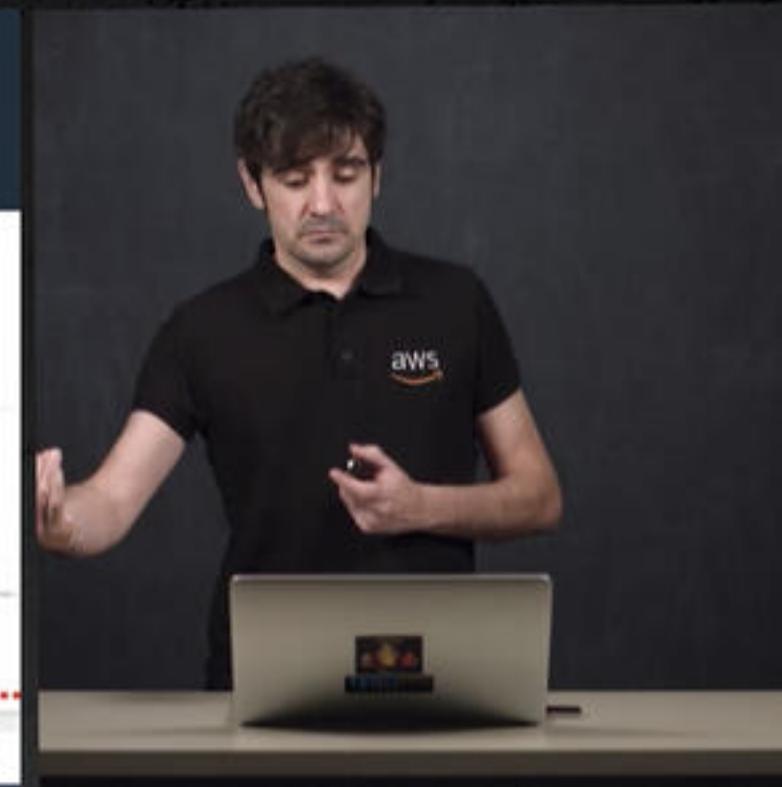
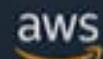
# Glue Jobs Optimization

Job Execution: Active Executors, Completed Stages & Maximum Needed Executors

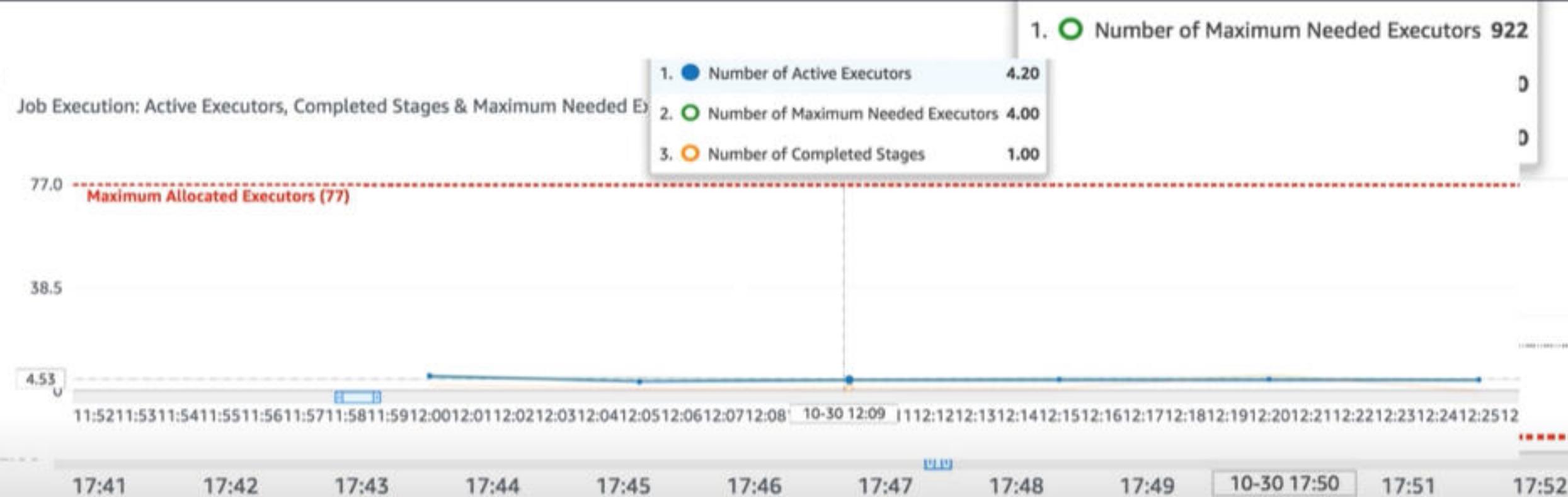
- 1. Number of Maximum Needed Executors 922
- 2. Number of Completed Stages 88.0
- 3. Number of Active Executors 57.0



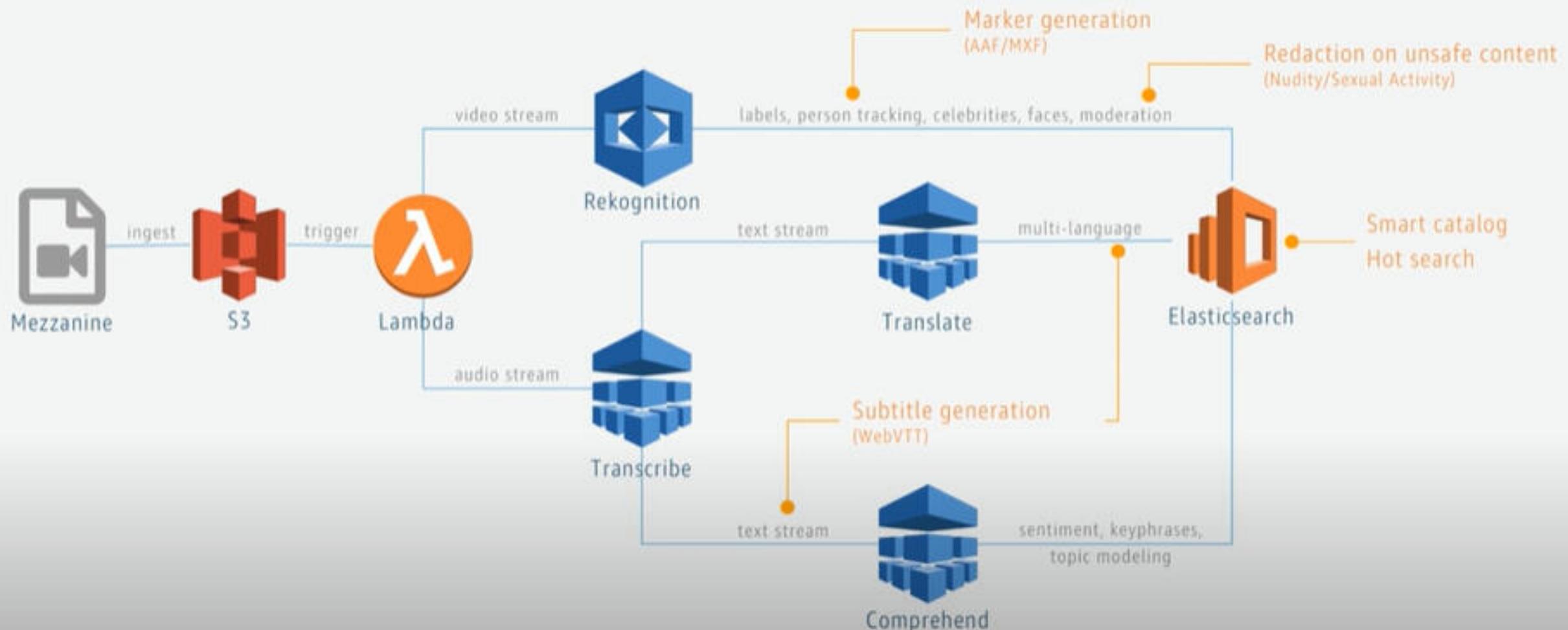
© 2018, Amazon Web Services, Inc. or its Affiliates. All rights reserved.



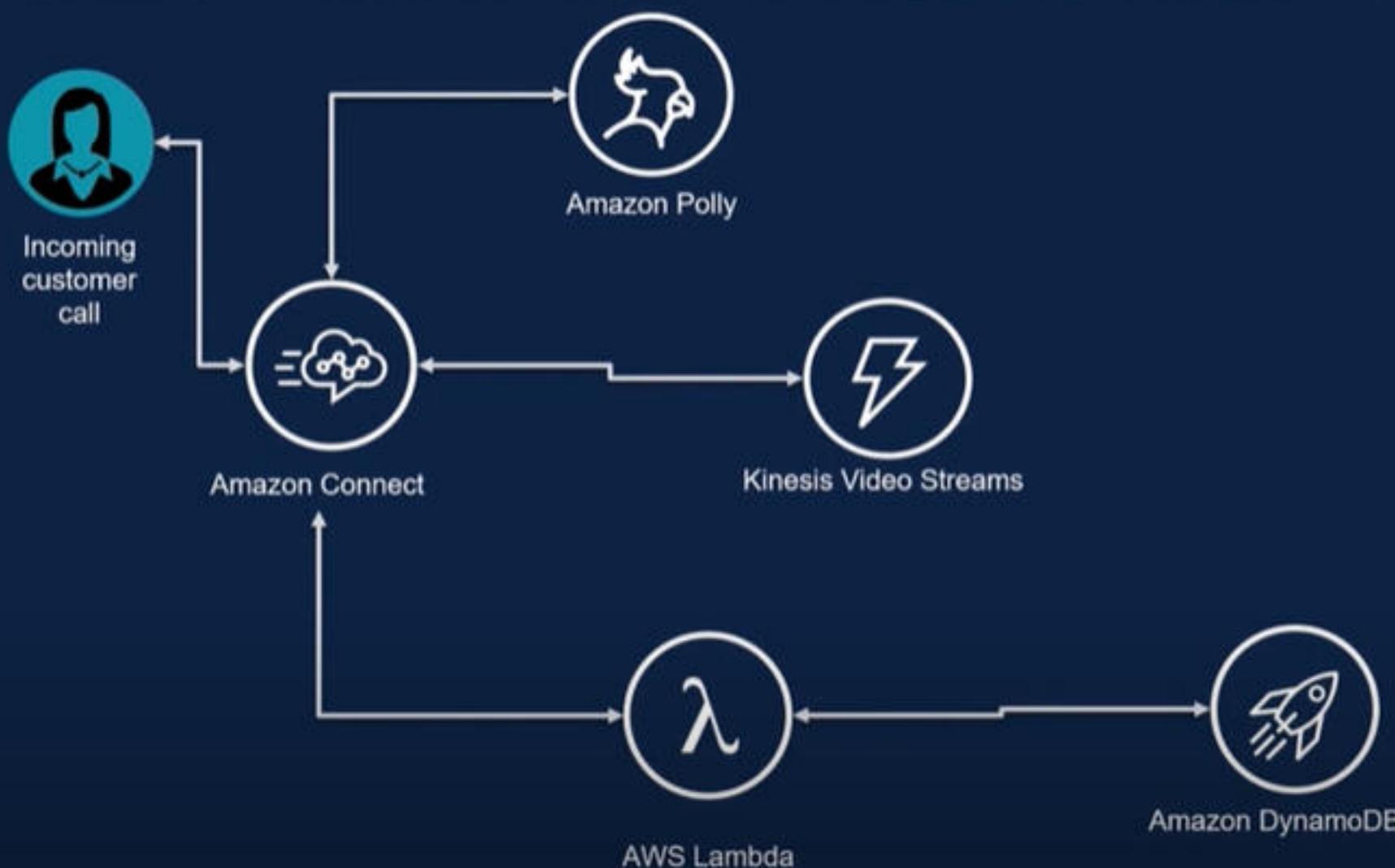
# Glue Jobs Optimization



# AI/ML use cases in M&E field



# Amazon Connect Realtime Customer Audio Transcription



details like the ARN, the

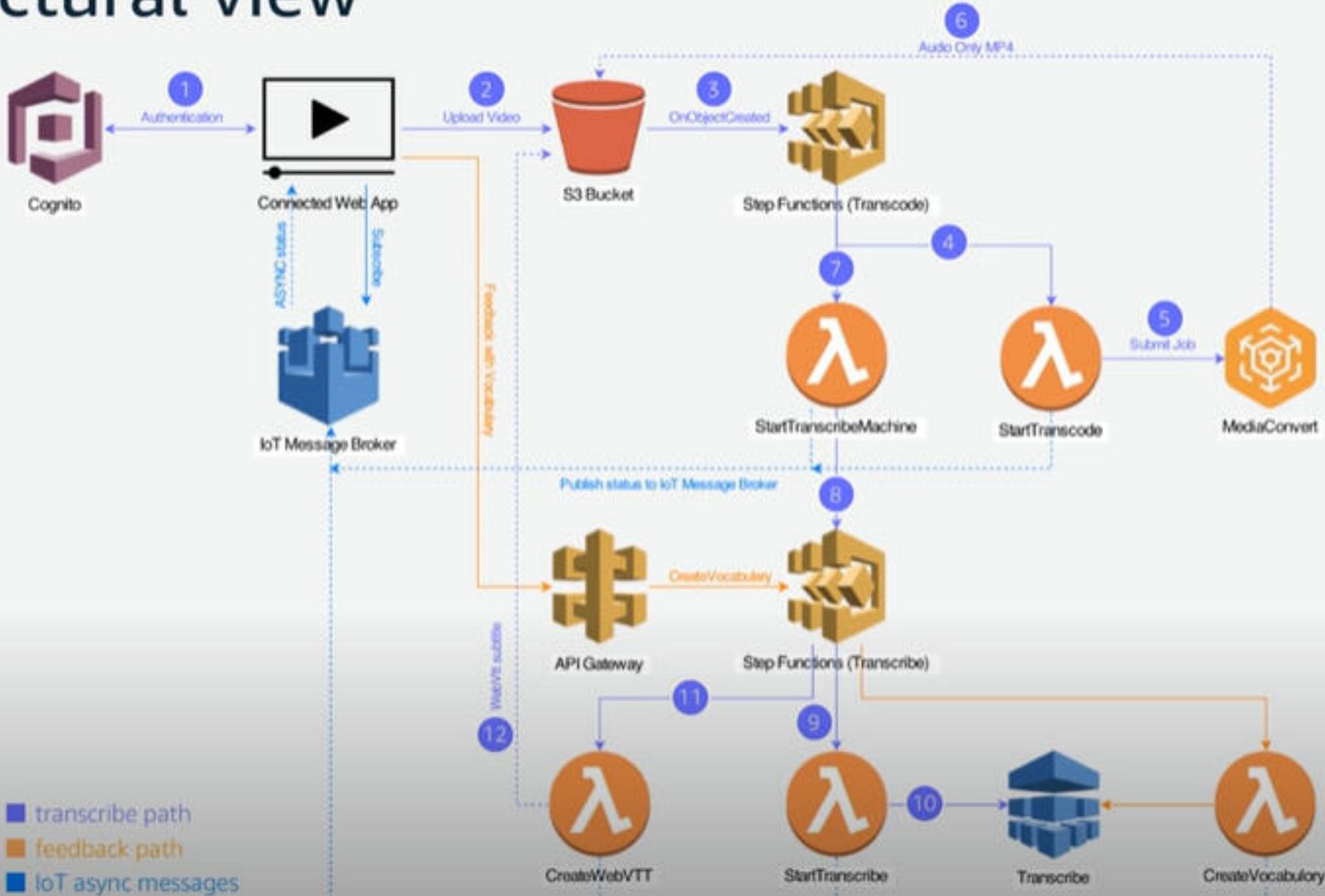


## MediaConvert: submit job (nodeJS)

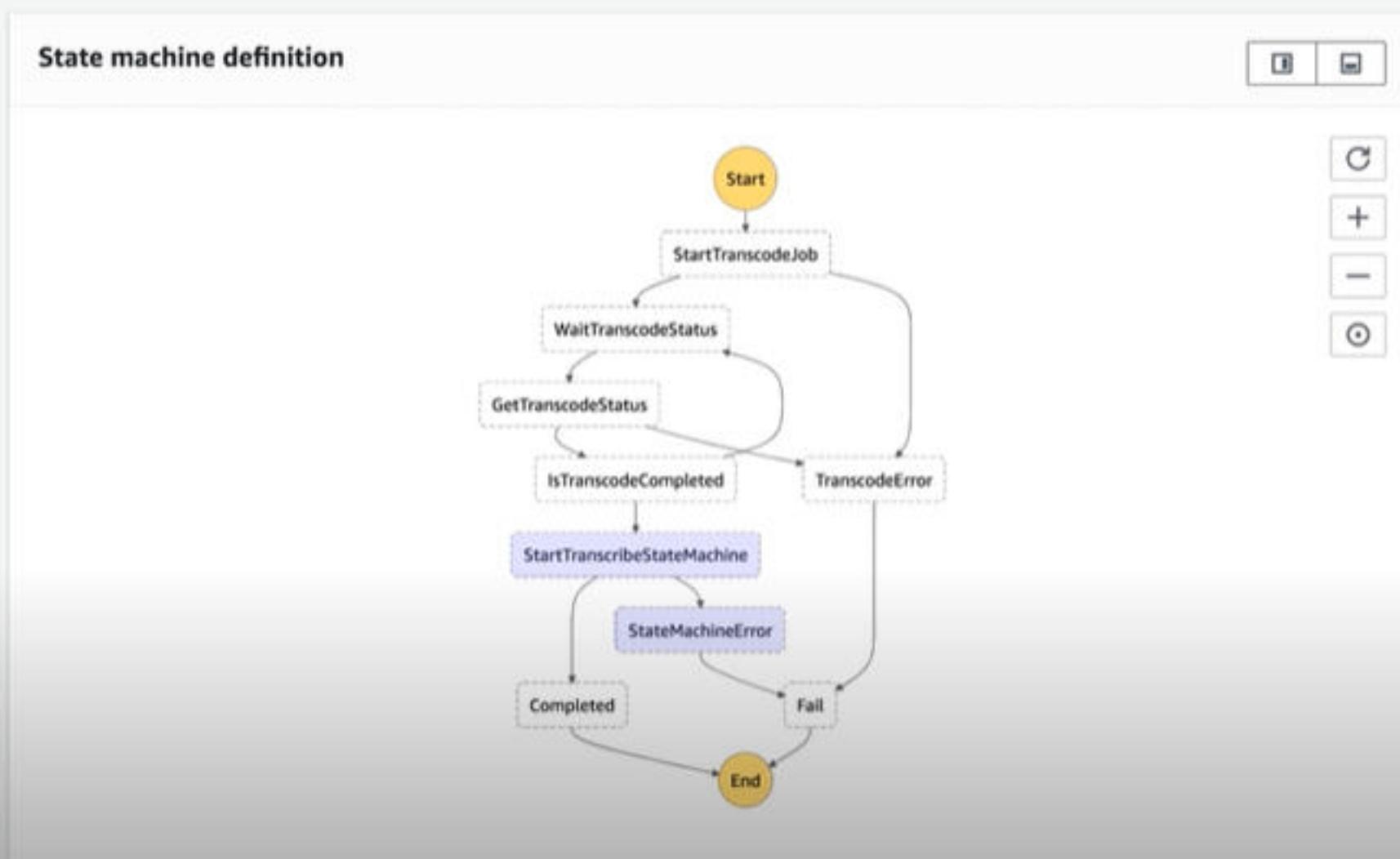
```
async function submitJob(template) {
  try {
    /* create a temp instance to find out the endpoint */
    const temporary = new AWS.MediaConvert({ apiVersion: '2017-08-29' });
    const data = await temporary.describeEndpoints().promise();
    const { Endpoints } = data;
    if (!Endpoints || !Endpoints.length) {
      throw new Error('no endpoint found!');
    }
    /* IMPORTANT: create a new instance and use the endpoint from describeEndpoint */
    const [{ Url: endpoint }] = Endpoints;
    const instance = new AWS.MediaConvert({
      apiVersion: '2017-08-29',
      endpoint,
    });
    const response = await instance.createJob(template).promise();
    return response;
  } catch (e) {
    throw e;
  }
}

const template = { /* Job JSON template */ };
const result = await submitJob(template);
```

# Architectural view



# Step Functions: transcoding state machine



# Transcribe: start transcribe service

```
const Basename = 'SOME_AUDIO_FILE';
_____

const params = {
    LanguageCode: 'en-US',
    MediaFormat: 'mp4',
    Media: { MediaFileUri: mediaFileUri },
    /* generate a unique transcribe name */
    TranscriptionJobName: `${Basename}-${new Date().toISOString().replace(/[-.:]/g, '')}`,
};

/* if vocabulary name exists, uses it */
if (Vocabularies) {
    params.Settings = { VocabularyName: Basename };
}
_____

const transcribe = new TranscribeService({ apiVersion: '2017-10-26' });
const response = await transcribe.startTranscriptionJob(params).promise();
```

# Transcribe: webvtt subtitle

WEBVTT

1  
00:00:00.040 --> 00:00:03.040  
Tell us which <c.five>ws</c> services  
2  
00:00:03.480 --> 00:00:06.480  
Absolutely. I love young <c.six>lumber</c> yard, which is a game  
3  
00:00:06.600 --> 00:00:09.600  
engine, that's really good, most like history, because  
4  
00:00:09.800 --> 00:00:12.800  
storage is fantastic, and, of course, you should really <c.five>lam</c>  
...

# IoT: mqtt message broker

Create a thing

Create IoT policy

Define a message topic

```
# Define your message topic.  
anyTopicName/status
```

# IoT: mqtt message broker

Create a thing

Create IoT policy

Define a message topic

Attach a policy to Cognito user

```
# Attach policy to Cognito user (identity)
aws iot attach-policy \
--policy-name anyThingPolicy \
--target "us-east-1:00000000-0000-0000-0000-000000000000"
```



# IoT: publisher

```
class IoTStatus {
  static async publish(endpoint, topic, payload) {
    try {
      const iotData = new AWS.IotData({ apiVersion: '2015-05-28', endpoint });
      const params = { topic, payload, qos: 0 };
      const response = await iotData.publish(params).promise();
    } catch (e) {
      e.message = `IoTStatus.publish: ${e.message}`;
    }
  }
}

/* Backend state machine to publish messages */
const endpoint = 'https://xxxxxxxxx-ats.iot.eu-west-2.amazonaws.com';
IoTStatus.publish(endpoint, 'anyTopicName/status', 'test message');
```

# IoT: subscriber

```
/* Download IoT JS SDK https://github.com/aws/aws-iot-device-sdk-js */
const subscriber = AWSIoTData.device({
  host: 'xxxxxxxxx-ats.iot.us-east-1.amazonaws.com',
  region: 'us-east-1',
  clientId: 'cognito-user',
  protocol: 'wss',
  debug: false,
  accessKeyId,
  secretKey,
  sessionToken,
});

subscriber.on('connect', () => {
  subscriber.subscribe('anyTopicName/status');
});

subscriber.on('message', (topic, payload) => {
  console.log(`Received: topic ${topic}: payload: ${payload.toString()}`);
  const message = JSON.parse(payload.toString());
  /* do something */
});
```

# MediaConvert: transcoding





## MediaConvert: IAM service role

```
# IAM role to allow MediaConvert to access S3/API resources
aws iam create-role \
--role-name MediaConvertServiceRole \
--assume-role-policy-document \
'{"Version":"2012-10-17", "Statement":[{"Sid":"","Effect":"Allow","Principal":{"Service":"mediaconvert.amazonaws.com"}, "Action":"sts:AssumeRole"}]}'

# Attach S3 policy
aws iam attach-role-policy \
--policy-arn arn:aws:iam::aws:policy/AmazonS3FullAccess \
--role-name MediaConvertServiceRole

# Attach API Gateway policy
aws iam attach-role-policy \
--policy-arn arn:aws:iam::aws:policy/AmazonAPIGatewayInvokeFullAccess \
--role-name MediaConvertServiceRole
```

# MediaConvert: define a job template

AWS Elemental MediaConvert > Jobs > Job details

## Job details

Summary

Duplicate

View JSON

Job

Inputs

Input 1

Output groups

mp4 - File group

Output 1

mp4 - File group

Output

Job settings

Settings

Show job JSON

### File group settings

Custom group name

mp4

Destination

c:\temp\el\transcoded\C-SPO\_R2D2\_Sesame\_Street\C-SPO\_R2D2\_Sesame\_Street

### Outputs

Output

Name modifier

Extension

Output 1

mp4

Audio only output (for Transcribe)

MP4 video output (for playback)



Sample JSON



# MediaConvert: submit job

```
# First, find out the per-region, per-account endpoint
aws mediaconvert describe-endpoints --region eu-west-1

{
    "Endpoints": [
        {
            "Url": "https://xxxxxxxx.mediaconvert.eu-west-1.amazonaws.com"
        }
    ]
}

# submit job to the endpoint
aws mediaconvert create-job \
--role arn:aws:iam::000000000000:role/MediaConvertServiceRole \
--settings file://sampleJob.json \
--endpoint https://xxxxxxxx.mediaconvert.eu-west-1.amazonaws.com \
--region eu-west-1

{
    JSON Job data...
}
```

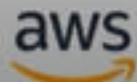
# MediaConvert: submit job

```
# First, find out the per-region, per-account endpoint
aws mediaconvert describe-endpoints --region eu-west-1

{
    "Endpoints": [
        {
            "Url": "https://xxxxxxxx.mediaconvert.eu-west-1.amazonaws.com"
        }
    ]
}

# submit job to the endpoint
aws mediaconvert create-job \
--role arn:aws:iam::000000000000:role/MediaConvertServiceRole \
--settings file://sampleJob.json \
--endpoint https://xxxxxxxx.mediaconvert.eu-west-1.amazonaws.com \
--region eu-west-1

{
    JSON Job data...
}
```



# MediaConvert: submit job (nodeJS)

```
async function submitJob(template) {
  try {
    /* create a temp instance to find out the endpoint */
    const temporary = new AWS.MediaConvert({ apiVersion: '2017-08-29' });
    const data = await temporary.describeEndpoints().promise();
    const { Endpoints } = data;
    if (!Endpoints || !Endpoints.length) {
      throw new Error('no endpoint found!');
    }
    /* IMPORTANT: create a new instance and use the endpoint from describeEndpoint */
    const [{ Url: endpoint }] = Endpoints;
    const instance = new AWS.MediaConvert({
      apiVersion: '2017-08-29',
      endpoint,
    });
    const response = await instance.createJob(template).promise();
    return response;
  } catch (e) {
    throw e;
  }
}

const template = /* Job JSON template */ ;
const result = await submitJob(template);
```

# A few gotcha...

1GB file size limit

*use MediaConvert, ETS, or ffmpeg to extract audio stream from source file*

2 hours duration limit

*split source file and post-process to stitch (modify timecode of) the transcription result*

working with acronym

*put a space in between acronym, AWS, EMR, and etc*

working with numbers

*pre/post-process to convert numbers before and after Transcribe*



# Evaluation of Accuracy: Likelihood



**Likelihood:** Likelihood is probability of some observed data—e.g. test, given a model.

$$P(\mathcal{E}_{test}; \theta) = \prod_{E \in \mathcal{E}_{test}} P(E; \theta)$$

## Example

$$|\mathcal{E}_{test}| = 10^6$$

E1 = I am from pittsburg,  $P(E; \theta) = 2.1 \times 10^{-6}$  \*

E2 = I study at a university  $P(E; \theta) = 1.1 \times 10^{-8}$

E3 = my mother is from utah,  $P(E; \theta) = 2 \times 10^{-9}$

$$P(\mathcal{E}_{test}; \theta) = \prod_{E \in \mathcal{E}_{test}} P(E; \theta) = 2.1 \times 10^{-6} \times 1.1 \times 10^{-8} \times 2.0 \times 10^{-9} = 4.62 \times 10^{-23}$$



08:44 / 42:33



# Evaluation of Accuracy: Log Likelihood



- *Theorem:*  $\log_b \prod_{i=1}^n x_i = \sum_{i=1}^n \log_b x_i$ ; *Example:*  $\log(x \cdot y) = \log(x) + \log(y)$
- We have seen probabilities of likelihood are very tiny. Multiplying these small values make the probabilities even smaller numbers. It is easier to work with log probabilities for measuring accuracy.

$$\log(P(\mathcal{E}_{test}; \theta)) = \log_b \prod_{E \in \mathcal{E}_{test}} P(E; \theta) = \sum_{E \in \mathcal{E}_{test}} \log(P(E; \theta)).$$

- **Example:**

$$|\mathcal{E}_{test}| = 10^6$$

$E_1 = \text{I am from pittsburg}$ ,  $P(E; \theta) = 2.1 \times 10^{-6} \therefore \log(P(E; \theta)) = \log(2.1) + (-6) = -5.67$ \*

$E_2 = \text{I study at a university}$   $P(E; \theta) = 1.1 \times 10^{-8} \therefore \log(P(E; \theta)) = \log(1.1) + (-8) = -7.96$

$E_3 = \text{my mother is from utah}$ ,  $P(E; \theta) = 2 \times 10^{-9} \therefore \log(P(E; \theta)) = \log(2.0) + (-9) = -8.70$

$$P(\mathcal{E}_{test}; \theta) = \prod_{E \in \mathcal{E}_{test}} P(E; \theta) = \sum_{E \in \mathcal{E}_{test}} \log(P(E; \theta)) = -5.67 - 7.96 - 8.70 = -22.33$$

# Evaluation of Accuracy: Perplexity



- It is common to divide log-likelihood by the number of words in the corpus. This makes cross-corpora measurement easier.

$$\text{length}(\mathcal{E}_{test}) = \sum_{E \in \mathcal{E}_{test}} |E|$$
$$ppl(\mathcal{E}_{test}; \theta) = e^{-(\log(P(\mathcal{E}_{test}; \theta)) / \text{length}(\mathcal{E}_{test}))}$$

- This accuracy indicator is called **perplexity**.

# NMT - Motivation

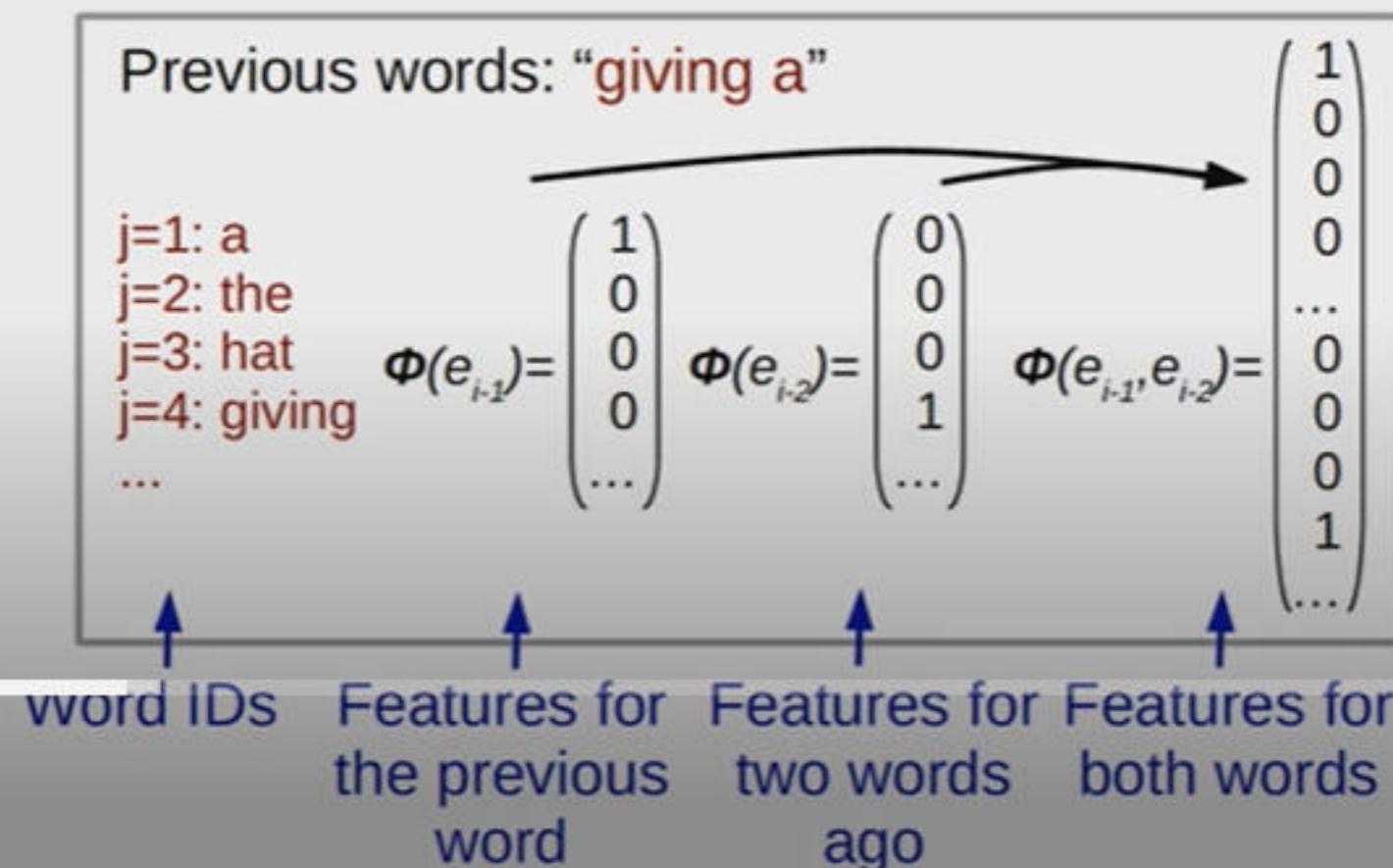


farmers eat	steak → <b>high</b> hay → <b>low</b>	cows eat	steak → <b>low</b> hay → <b>high</b>
farmers grow	steak → <b>low</b> hay → <b>high</b>	cows grow	steak → <b>low</b> hay → <b>low</b>

- In separate contexts  $e_{t-2} = \text{"farmers"}$  is compatible with  $e_t = \text{"hay"}$  and  $e_{t-1} = \text{"eats"}$ .
- If we are using feature sets depending on  $e_{t-1}$  and  $e_{t-2}$ , then a sentence like "farmers eat hay" cannot be ruled out, even though it is an unnatural sentence or perhaps an unnatural farmer.

# Step1: Word Embedding

- In Neural Language Models we represent words as word embedding, **a vector representation of words.**
- Then we concatenate these words to the length of the context (3 for trigrams).



# Step3: Score and Probability

- Next, we calculate the score vector for each word. This is done by performing an affine\* transform of the hidden vector  $h$  with a weight matrix and adding a bias vector.
- Finally, we get a probability estimate  $p$  by running the calculated scores through a softmax function.

$$\mathbf{m} = \text{concat}(M_{\cdot, e_{t-2}}, M_{\cdot, e_{t-1}})$$

$$\mathbf{h} = \tanh(W_{mh}\mathbf{m} + \mathbf{b}_h)$$

$$\mathbf{s} = W_{hs}\mathbf{h} + \mathbf{b}_s$$

$$\mathbf{p} = \text{softmax}(\mathbf{s})$$

\* An affine transformation is any transformation that preserves collinearity (i.e., all points lying on a line initially still lie on a line after transformation) and ratios of distances (e.g., the midpoint of a line segment remains the midpoint after transformation). Ref:

# Benefits of Neural Models Contd.



- **Ability to skip previous words:** n-gram models refer to all sequential words depending the length of the context. Neural Models can skip words.
- **Simplistic Models:** SMT models require complex data-structures, in production often using 80 GBs or more of heap space. High quality NMT production models are generally less than 1 GB.

# Dependencies: Reference and Context Cont.



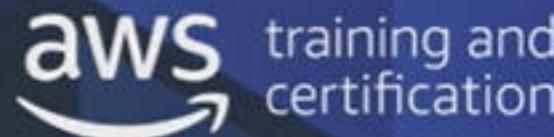
- **The Spanish King** officially abdicated in **favour** of his **son, Felipe**. **Felipe** will be confirmed tomorrow as the new Spanish **king**. As the new **king, he** will be the sovereign.

Reference

Context

Reference & Context

# Dependencies: Reference and Context Cont.



- **The Spanish King** officially abdicated in **favour** of his **son, Felipe**. **Felipe** will be confirmed tomorrow as the new Spanish **king**. As the new **king, he** will be the sovereign.
- Their duo was known as "Buck and Bubbles". Bubbles tapped along as Buck **played on the stride piano and sang** until Bubbles was totally tapped out.

Reference

Context

Reference & Context



18:36 / 42:33



# Problems of Translation



To perform translation we need to solve the following problems:

**Modeling:** Decide what our model  $P(E|F; \theta)$  will look like. What parameters will it have, and how will the parameters specify a probability distribution?

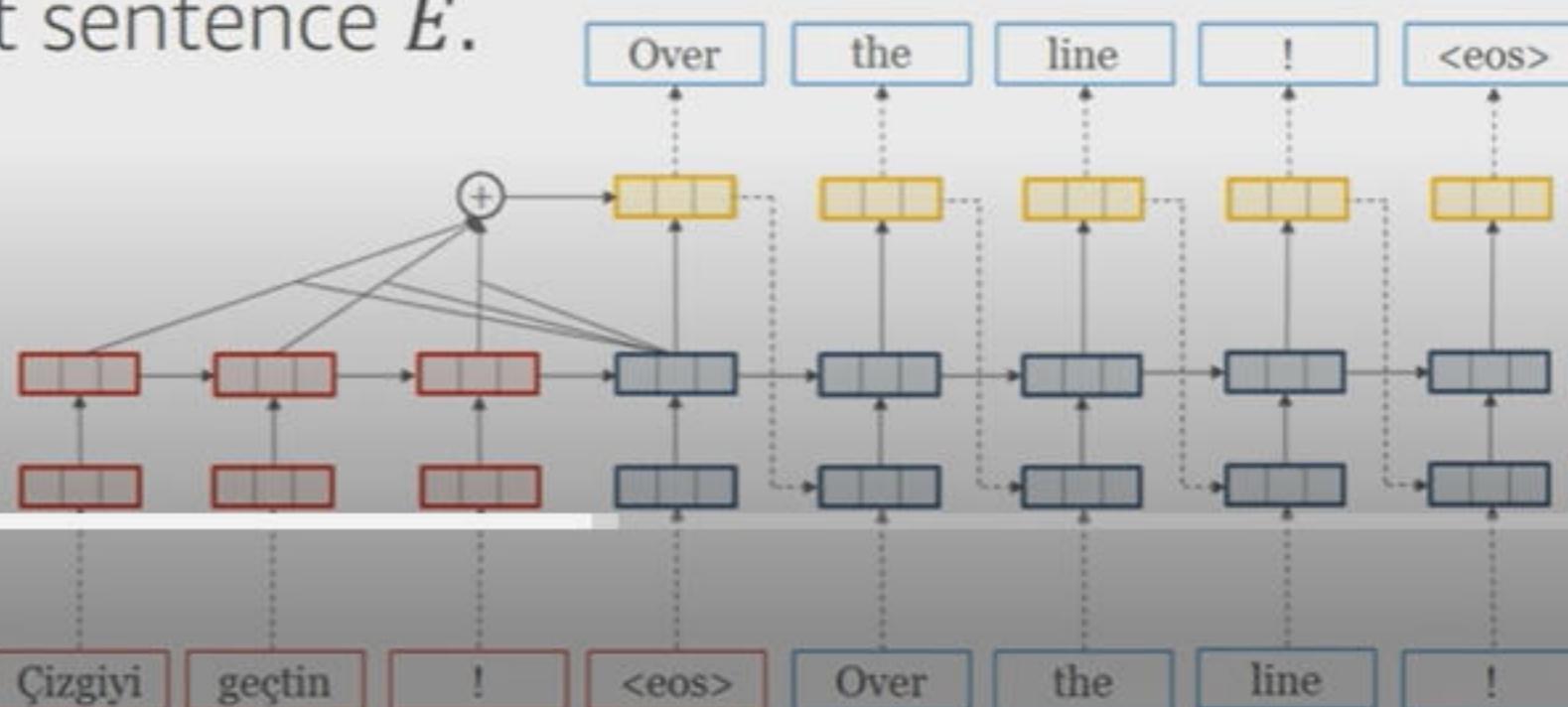
**Learning:** We need a method to learn appropriate values for parameters from training data.

**Search:** We need to find the most probable sentence (solving "argmax"). This process of searching for the best hypothesis is often called decoding.

# Inferring on Sequence from Another

We still intend to calculate  $P(E; \theta)$ , but before doing so we would like to calculate the initial state of the target model based on another RNN over the source sentence  $F$ .

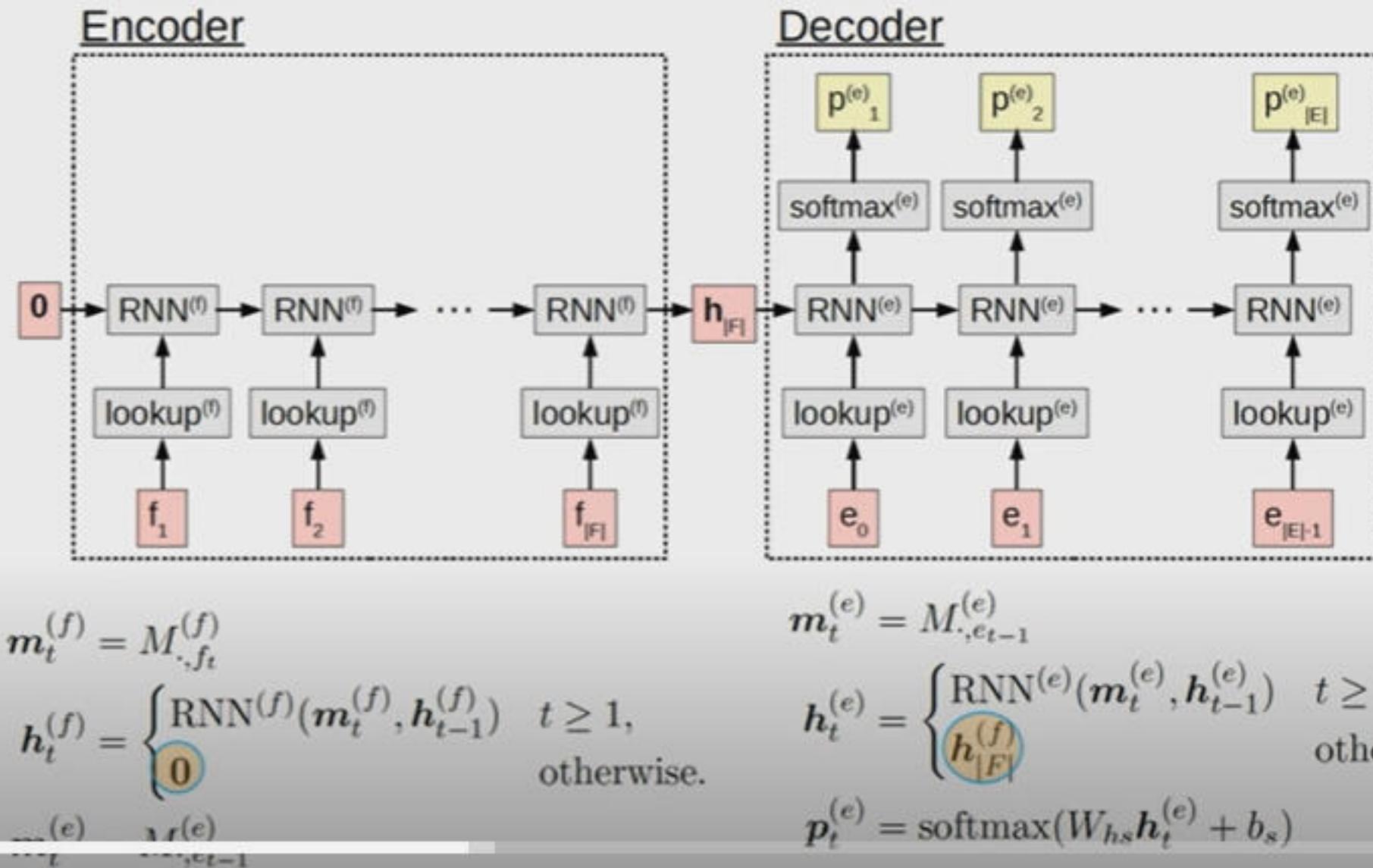
First RNN "encodes" source sentence  $F$ , and second RNN "decodes" the initial state into target sentence  $E$ .



# Encoder-Decoder Architecture



A model based on two LSTM networks, one that takes a sentence and captures it into an output layer that is to be used as input for another LSTM that performs translation



# Generating Output

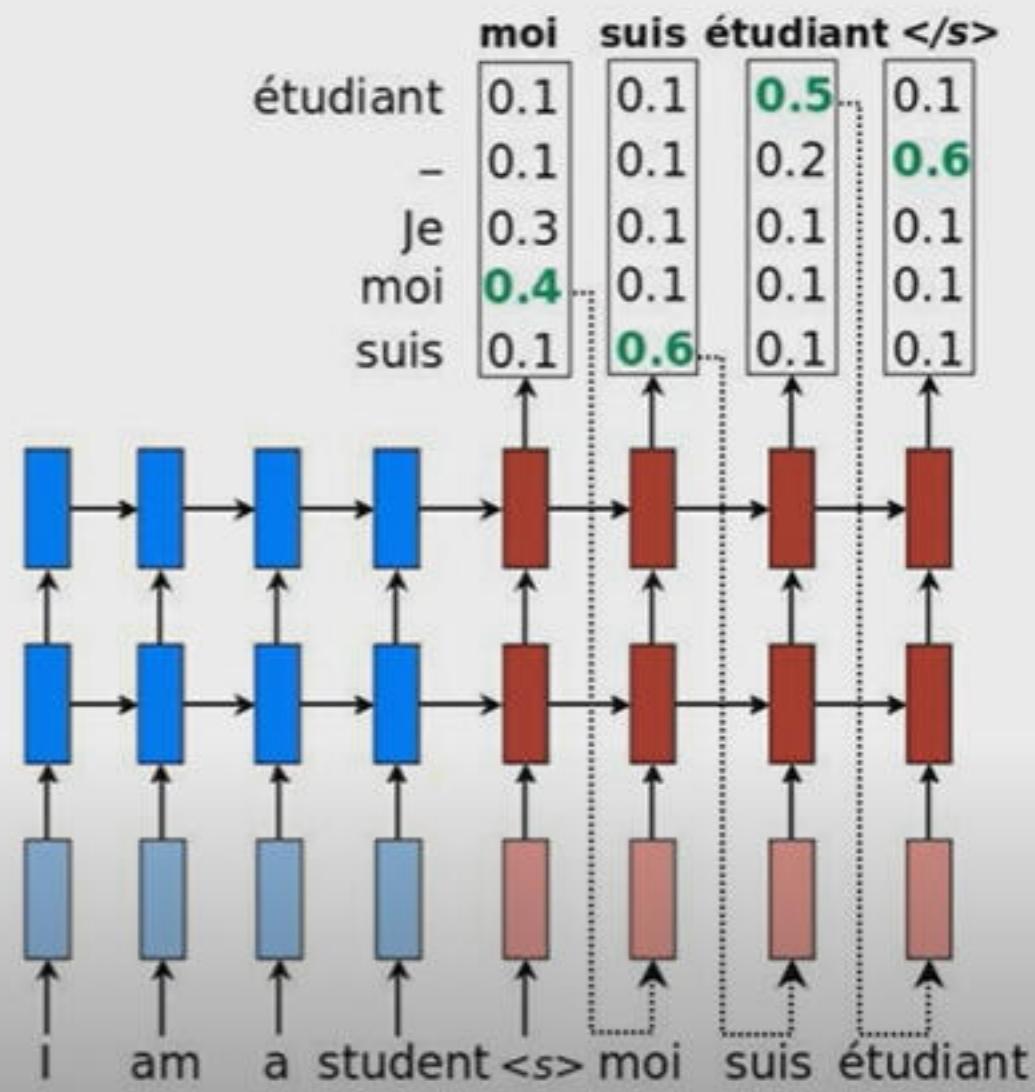
- So far we have seen how to make a probability model. In order to perform translation, we need to generate output. Generally this is performed using search in **parallel corpora**, the corpora of target language.
- Search is performed based on several Criteria:
  - **Random Sampling:** Randomly selecting an output  $E$  from a probability model.  $\hat{E} \sim P(E|F)$
  - **Greedy 1-Best Search:** Finding  $E$  that maximizes  $P(E|F)$ ;  $\hat{E} = \operatorname{argmax}_E P(E|F)$
  - **Beam Search:** Finding the n options with the highest probability  $P(E|F)$

# Greedy (Best) Search

- This is very similar to random sampling with one difference: The goal is to search for the word in target corpora that maximizes probability.

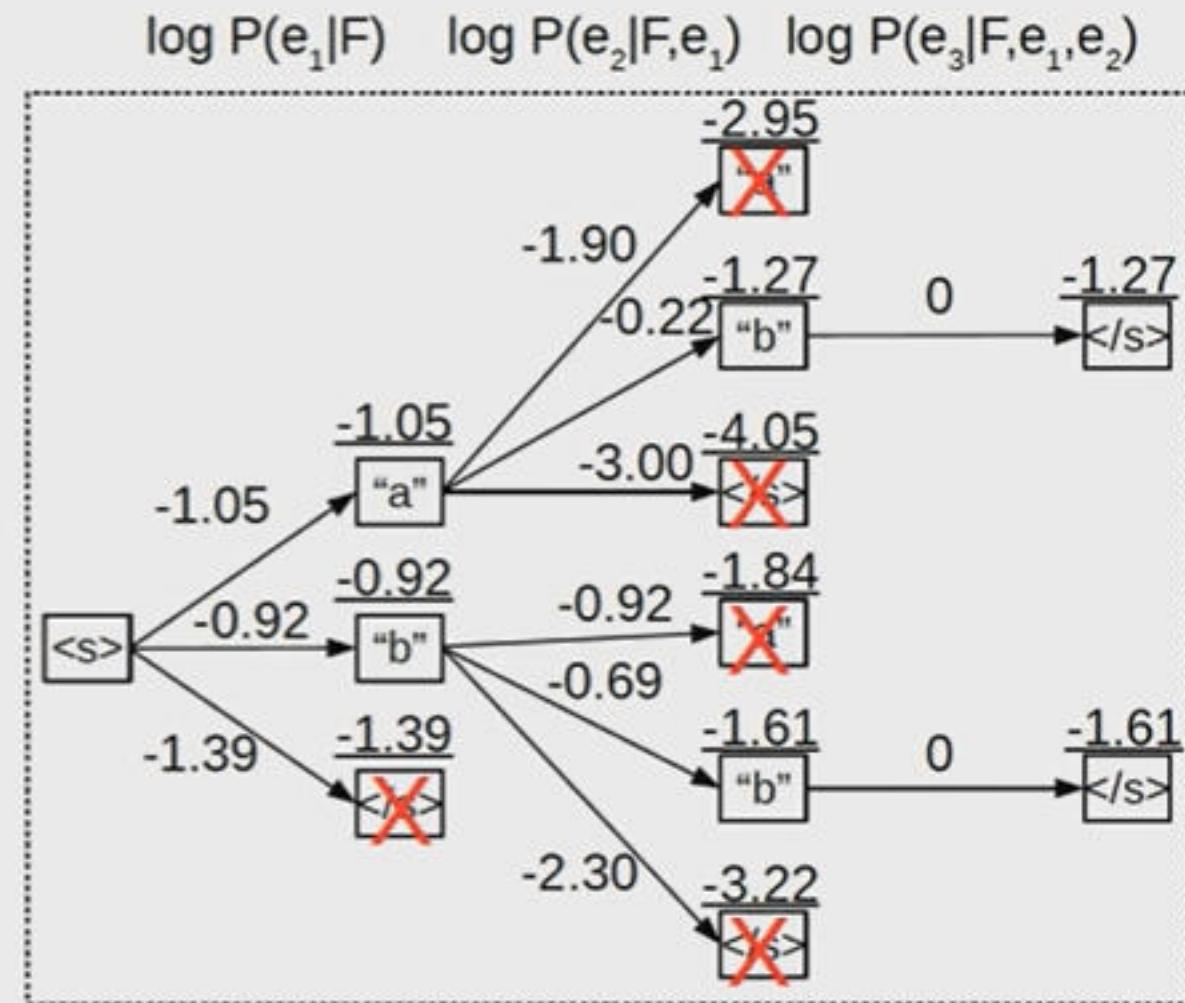
$$\hat{e}_t = \operatorname{argmax}_i \log(p_{t,i}^{(e)})$$

- Greedy search is not guaranteed to find the translation with highest probability, but is efficient in terms of memory and computation.



# Beam Search

- Beam Search is a search method that searches for  $b$  best hypothesis at each time step.  $b$  is called width of the beam.
- Larger beam size has a strong length bias.
- $b$  is a hyperparameter that should be selected to maximize accuracy.
- Not very efficient.



# Limitations

- Long distance dependencies.
- Storing variable length sentences in fixed size vectors.

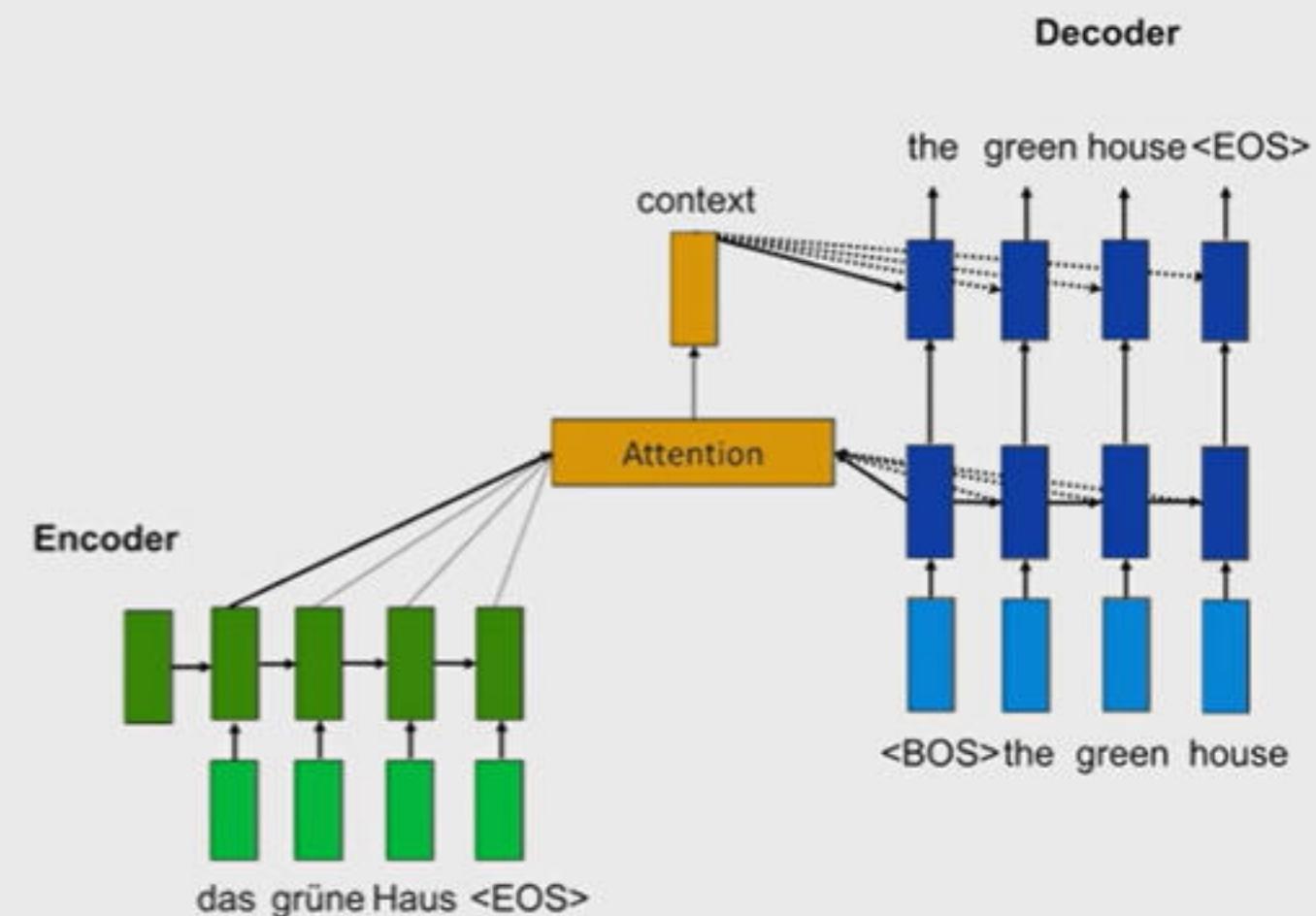
He purchased a car, so that he **could** drive to work and not spend his time on train.



Er kaufte ein Auto, damit er zur Arbeit fahren und seine Zeit nicht im Zug verbringen **konnte**

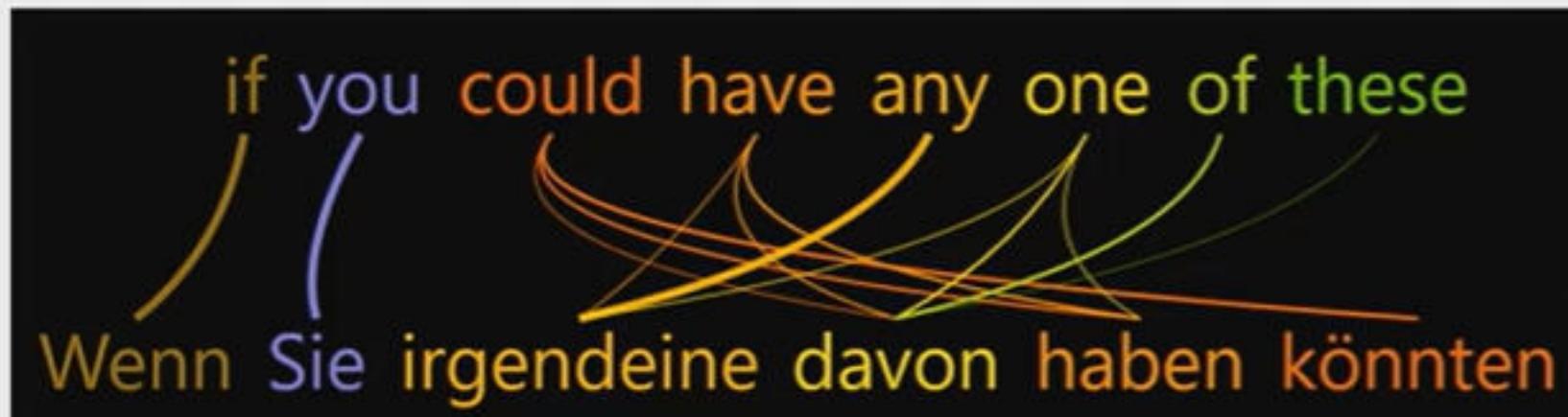
# Attention

- The attention model comes between the encoder and the decoder and helps the decoder to pick only the encoded inputs that are important for each step of the decoding process.
- For each encoded input from the encoder RNN, the attention mechanism calculates its importance.

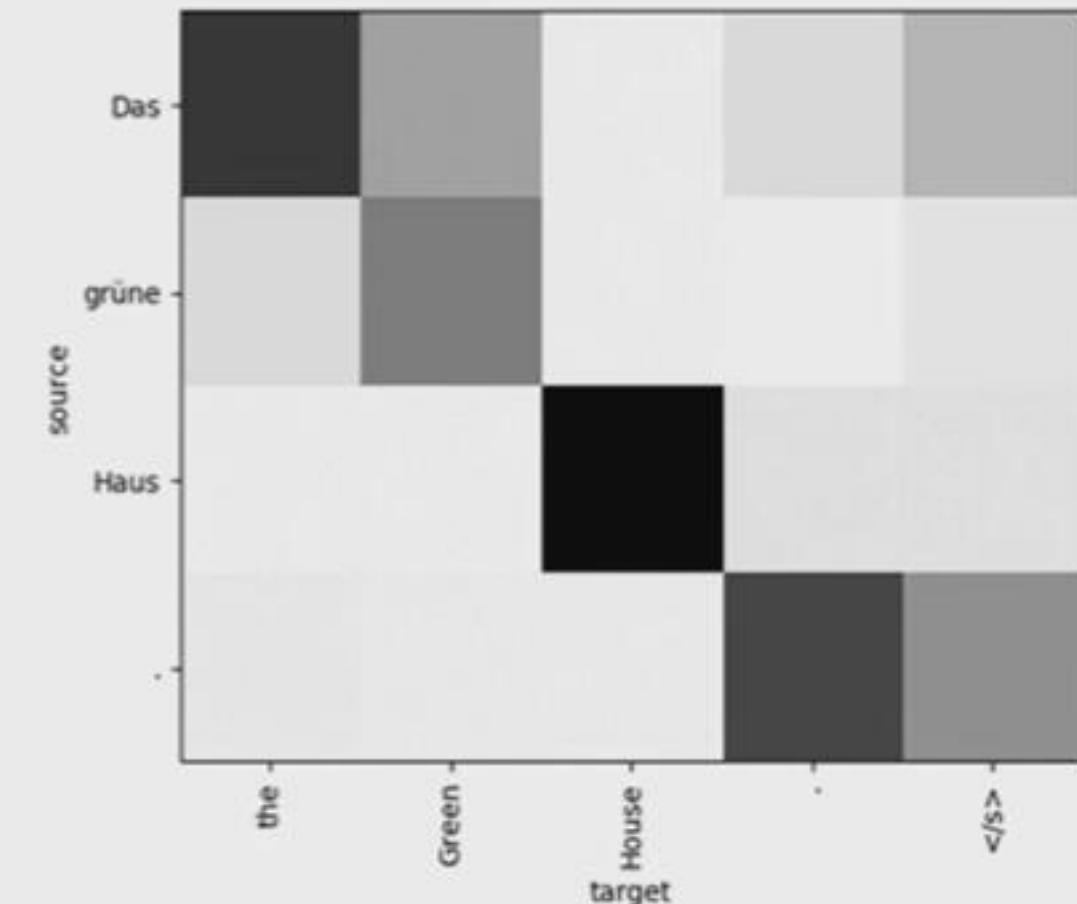


<https://aws.amazon.com/blogs/ai/train-neural-machine-translation-models-with-sockeye/>

# Visualization of Importance



<https://aws.amazon.com/blogs/ai/amazon-at-wmt-improving-machine-translation-with-user-feedback/>

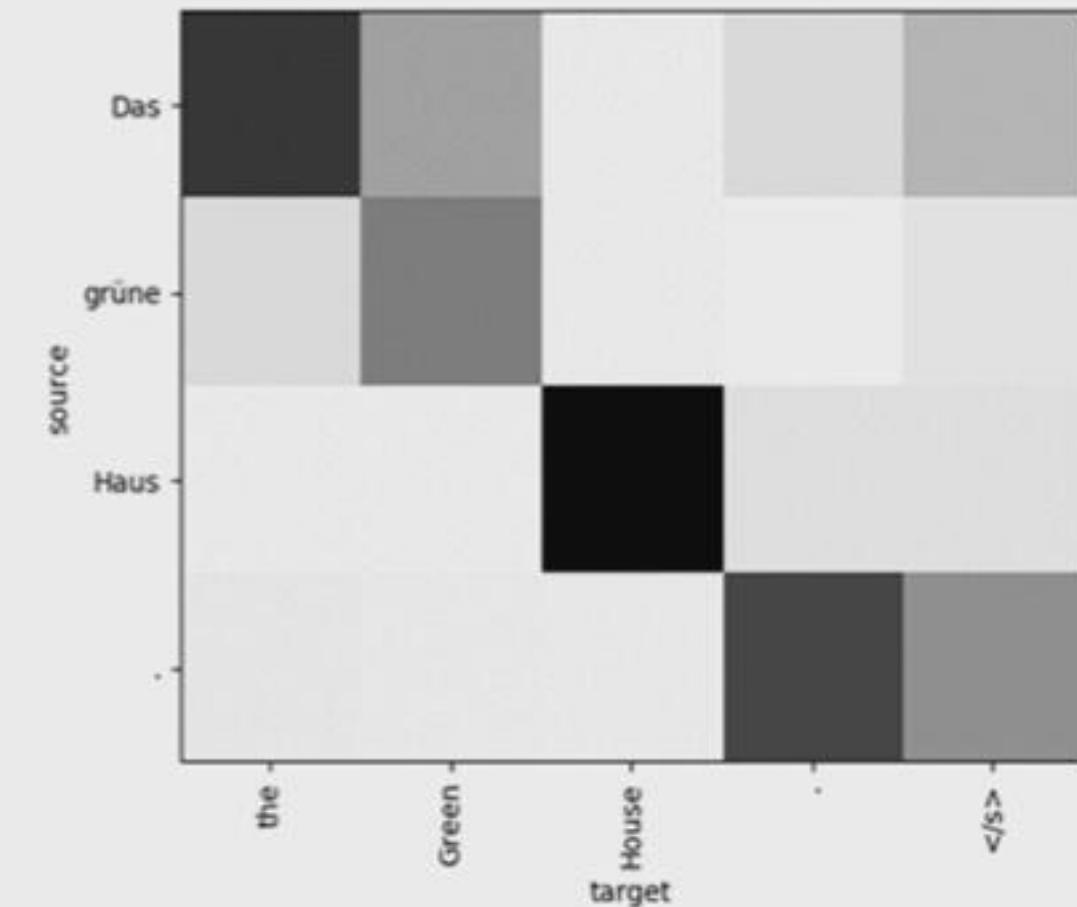


<https://aws.amazon.com/blogs/ai/train-neural-machine-translation-models-with-sockeye/>

# Visualization of Importance



<https://aws.amazon.com/blogs/ai/amazon-at-wmt-improving-machine-translation-with-user-feedback/>



<https://aws.amazon.com/blogs/ai/train-neural-machine-translation-models-with-sockeye/>

# Attention Mechanism



The attention computation happens at every decoder time step. It consists of the following stages:

1. The current target hidden state is compared with all source states to derive *attention weights*.
2. Based on the attention weights we compute a *context vector* as the weighted average of the source states.
3. Combine the context vector with the current target hidden state to yield the final *attention vector*.
4. The attention vector is fed as an input to the next time step (*input feeding*). The first three steps can be summarized by the following equations:

# Attention Mechanism Contd.

$$\alpha_{ts} = \frac{\exp(\text{score}(\mathbf{h}_t, \bar{\mathbf{h}}_s))}{\sum_{s'=1}^S \exp(\text{score}(\mathbf{h}_t, \bar{\mathbf{h}}_{s'}))}$$

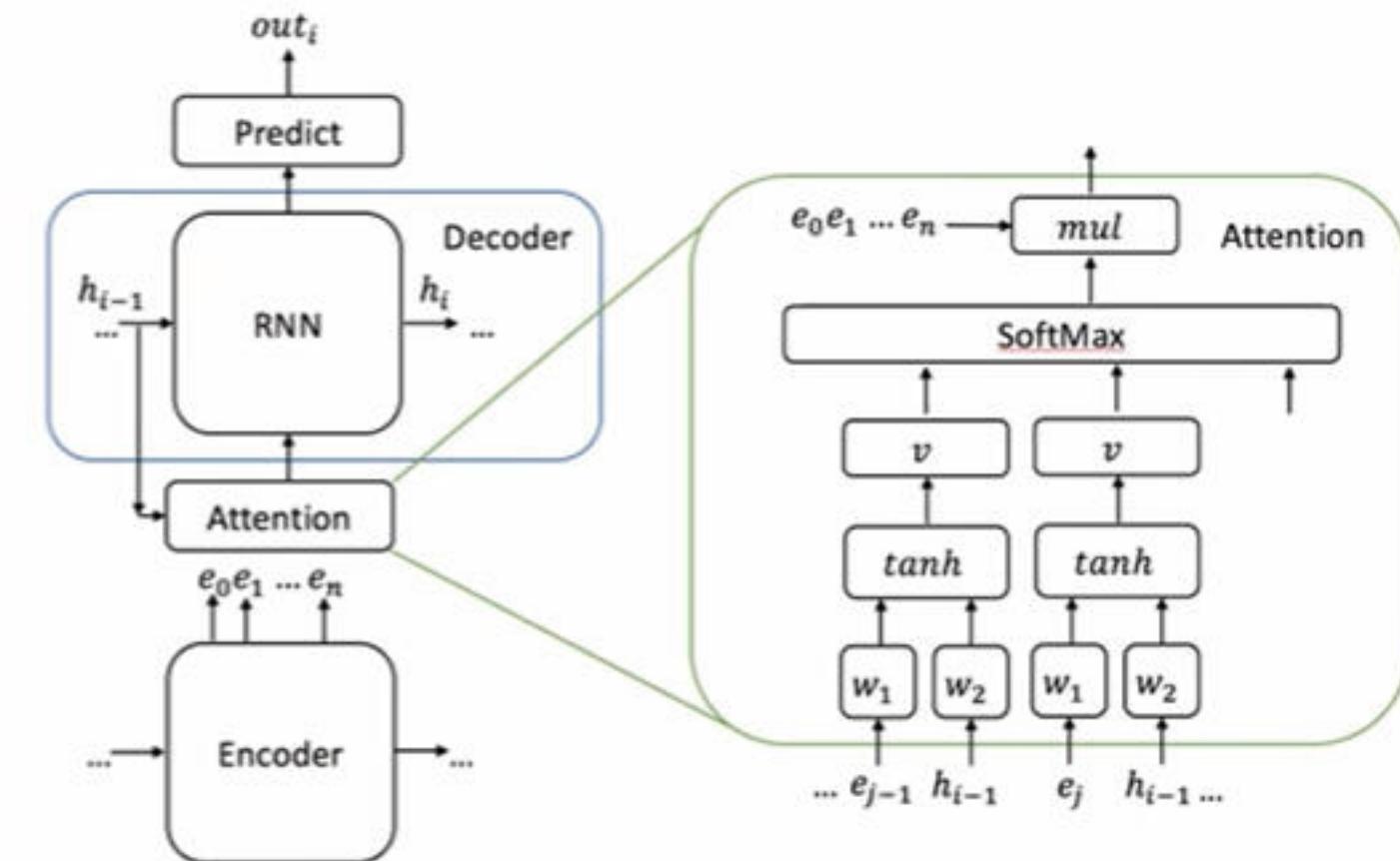
$$\mathbf{c}_t = \sum_s \alpha_{ts} \bar{\mathbf{h}}_s$$

$$\mathbf{a}_t = f(\mathbf{c}_t, \mathbf{h}_t) = \tanh(\mathbf{W}_c[\mathbf{c}_t; \mathbf{h}_t])$$

$\alpha_{ts}$ : attention weights

$c_t$ : context vector

$a_t$  = attention vector



<https://talbaumel.github.io/attention/>

# Calculating Attention Score



Method	Pros	Cons	Formula
Dot Product	<ul style="list-style-type: none"><li>- Simple</li><li>- No additional Parameters</li></ul>	<ul style="list-style-type: none"><li>- Input and output encoding in the same space</li></ul>	$\text{attn\_score}\left(h_j^{(f)}, h_t^{(e)}\right) := h_j^{(f)T} h_t^{(e)}$
Bilinear Functions	<ul style="list-style-type: none"><li>- Input and output of different sized</li></ul>	<ul style="list-style-type: none"><li>- More expensive than <math>\otimes</math></li></ul>	$\text{attn\_score}\left(h_j^{(f)}, h_t^{(e)}\right) := h_j^{(f)T} W_a h_t^{(e)}$ <i>W<sub>a</sub> is a linear transformation</i>
MLP	<ul style="list-style-type: none"><li>- More flexible than <math>\otimes</math></li><li>- Fewer parameters than BLF</li></ul>		$\text{attn\_score}\left(h_j^{(f)}, h_t^{(e)}\right) :=$ $w_{a2}^T \tanh\left(W_{a1}[h_t^{(e)}; h_j^{(f)}]\right)$ <i>W<sub>a1</sub> and w<sub>a2</sub> are weight matrix and vector of the first and second layers of the MLP respectively</i>

# Translation Performance



- **BLEU:** The BLEU score measures similarity to human translation. It measures how many words overlap in a *given translation* when compared to a *reference translation*. BLEU scores range from 0-100, the higher the score, the more the translation correlates to a human translation.
- **TER:** The TER score measures the amount of editing that a translator would have to perform to change a translation so it exactly matches a reference translation. By repeating this analysis on a large number of sample translations, it is possible to estimate the post-editing effort required for a project. TER scores also range from 0-100.

# Translation Performance Contd.



criterion	loss	BLEU	TER	NIST
MLE	N/A	30.48	60.85	8.26
MRT	-sBLEU	<b>36.71</b>	<b>53.48</b>	<b>8.90</b>
	sTER	30.14	53.83	6.02
	-sNIST	32.32	56.85	<b>8.90</b>

# Sockeye Encode-Decoder Framework



- Sockeye is a sequence-to-sequence framework with attention for Neural Machine Translation based on Apache MXNet Incubating. It implements the well-known encoder-decoder architecture with attention.
- Github repo: <https://github.com/awslabs/sockeye>
- Tutorials: <https://github.com/awslabs/sockeye/tree/master/tutorials>

# Examples

- `python3 -m sockeye.train -s train.source \ -t train.target \ -vs dev.source \ -vt dev.target \ --num-embed 32 \ --rnn-num-hidden 64 \ --attention-type dot \ --use-cpu \ --metrics perplexity accuracy \ --max-num-checkpoint-not-improved 3 \ -o seqcopy_model`
- `python -m sockeye.translate -m seqcopy_model`

# Sockeye Features



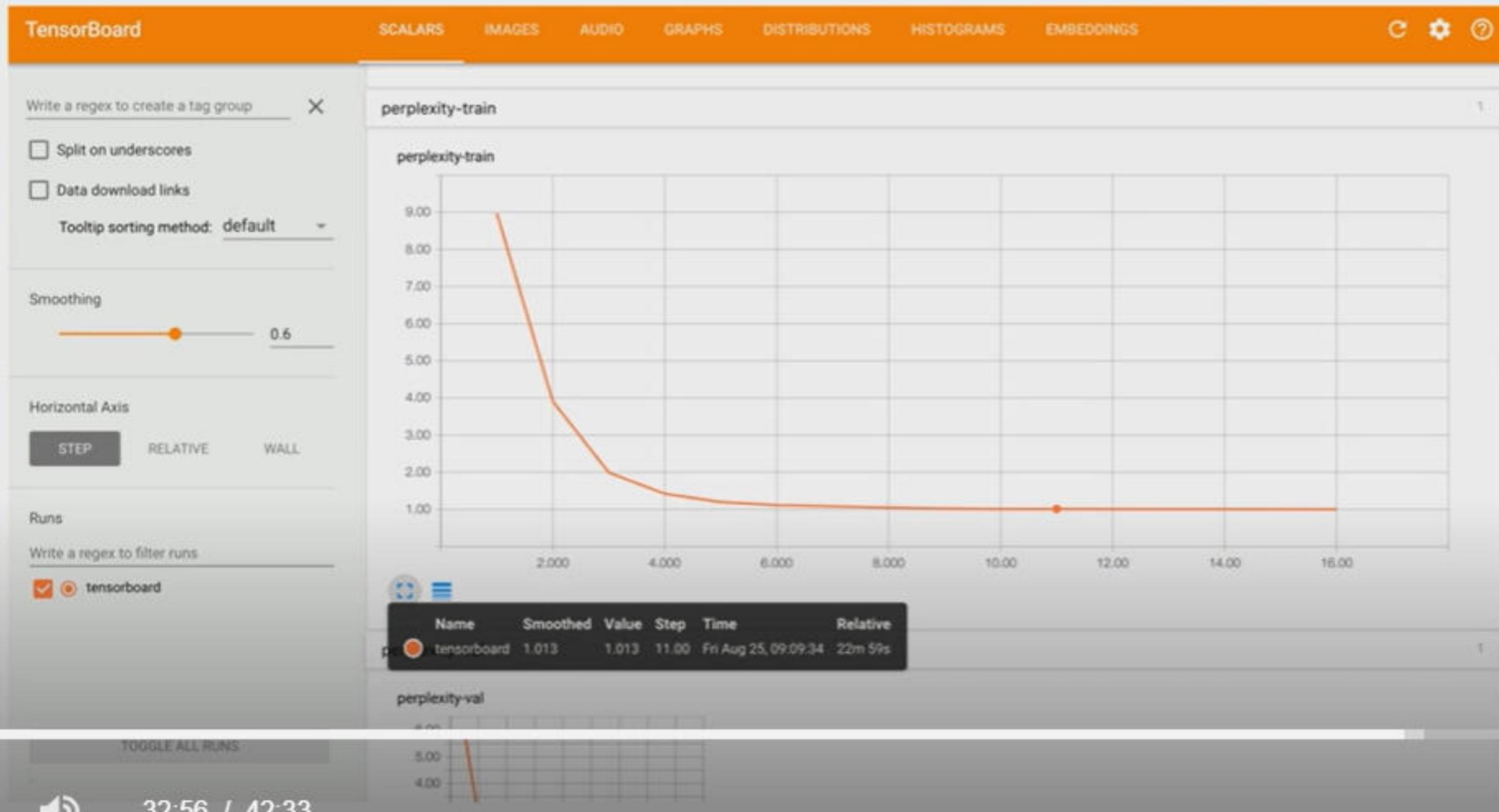
- Beam search inference
- Easy ensembling of multiple models
- Residual connections between RNN layers ([Wu et al., 2016](#)) [Deep LSTM with parallelism]
- Lexical biasing of output layer predictions ([Arthur et al., 2016](#)) [Low Frequency Words]
- Modeling coverage ([Tu et al., 2016](#)) [Keeping attention history to reduce over and under translation]
- Context gating ([Tu et al., 2017](#)) [Improving adequacy of translation by controlling ratios of source and target context]

# Sockeye Features Contd.

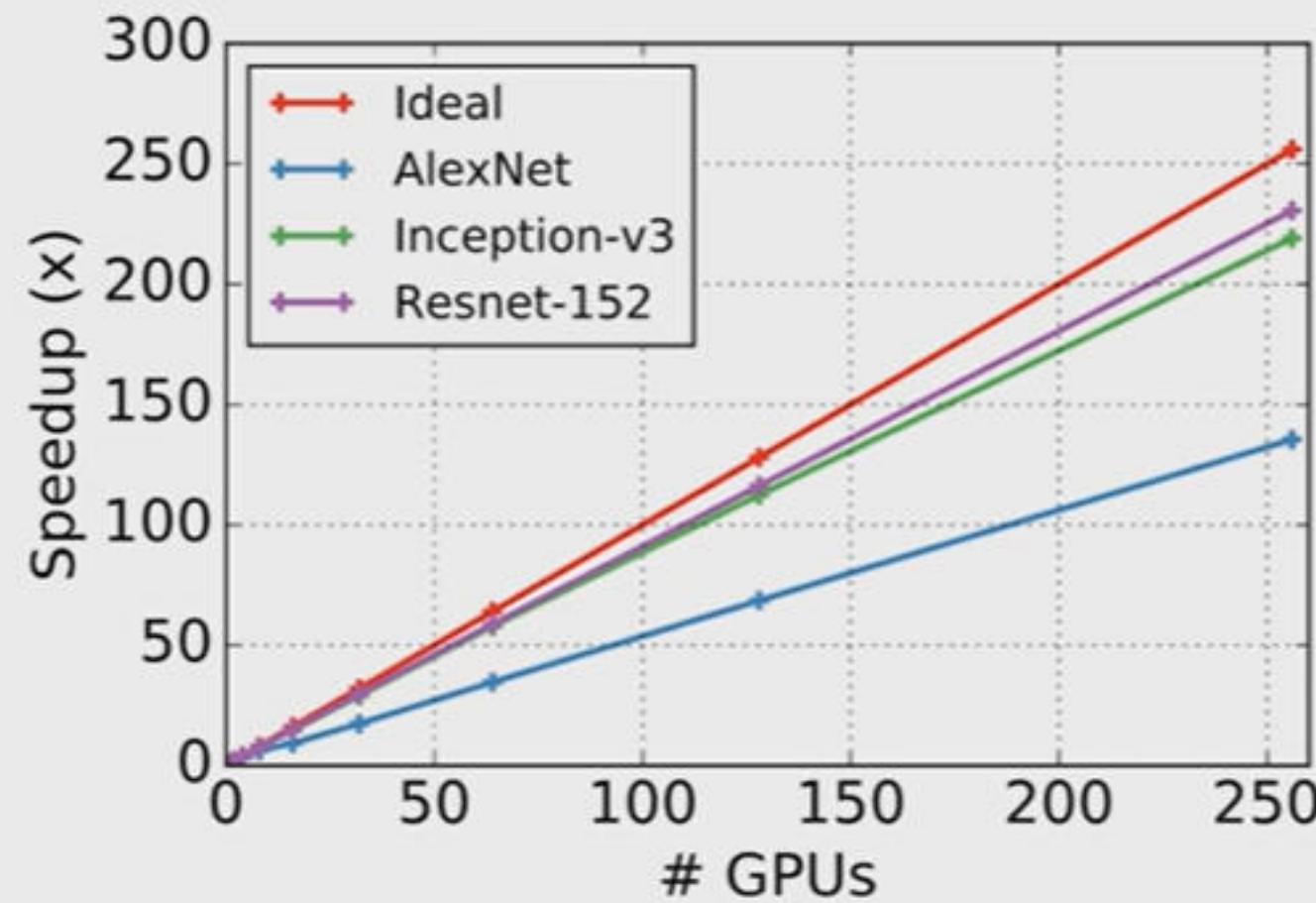


- Cross-entropy label smoothing (e.g., [Pereyra et al., 2017](#))
- Layer normalization ([Ba et al, 2016](#)) [improving training time]
- Multiple supported attention mechanisms [dot, mlp, bilinear, multihead-dot, encoder last state, location]
- Multiple model architectures (Encoder-Decoder [Wu et al., 2016](#), Convolutional [Gehring et al, 2017](#), Transformer [Vaswani et al, 2017](#))

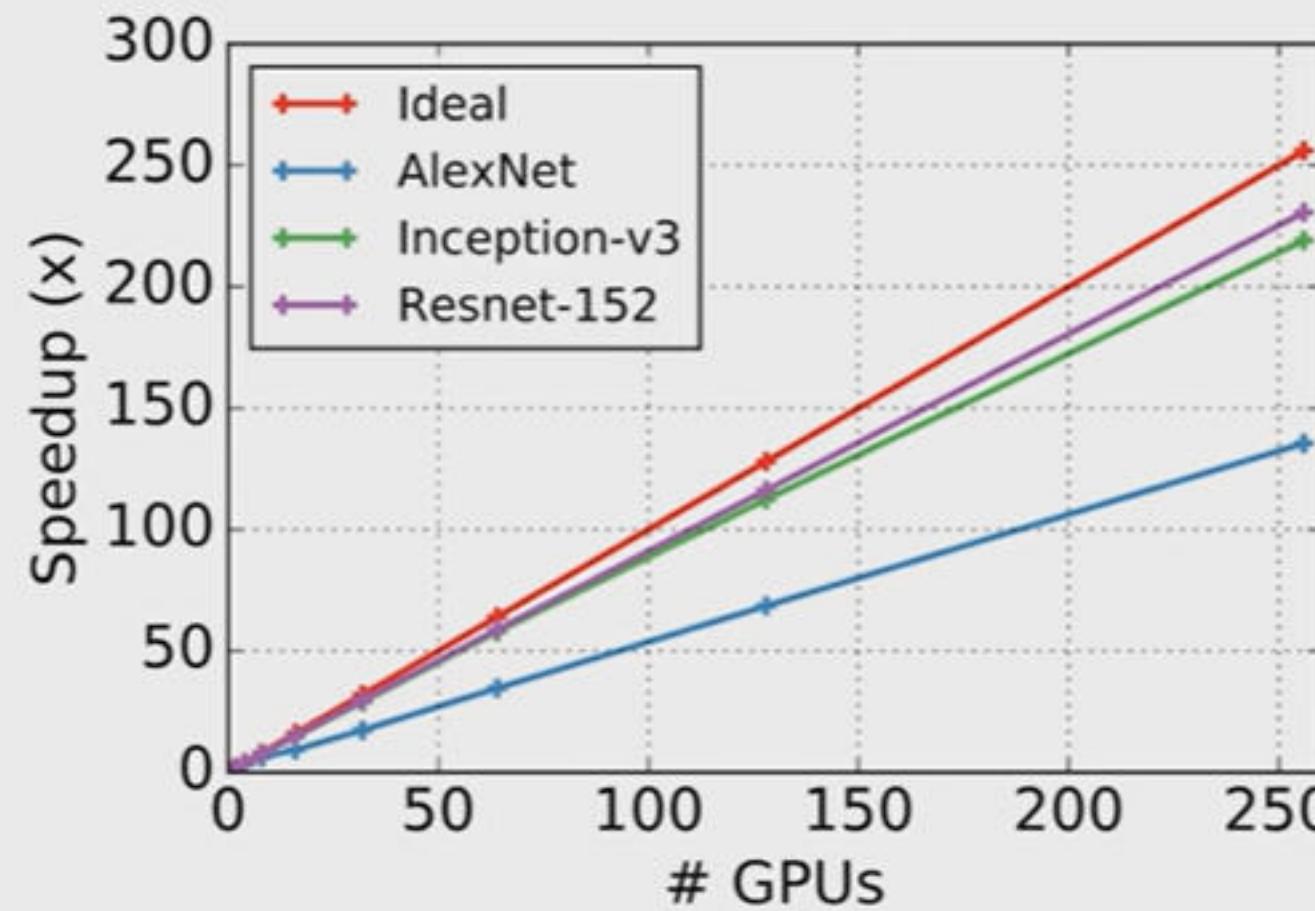
# Metrics and Visualization



- Sockeye uses MXNet as an underlying deep learning framework.
- MXNet benchmark shows near linear performance across multiple machines.



- Sockeye uses MXNet as an underlying deep learning framework.
- MXNet benchmark shows near linear performance across multiple machines.



# Multi-GPU Training

Training can be carried out on multiple GPUs by either specifying multiple GPU device ids: --device-ids 0 1 2 3, or specifying the number GPUs required: --device-ids -n attempts to acquire n GPUs.

Every 2.0s: nvidia-smi						
Wed Oct 11 22:28:25 2017						
NVIDIA-SMI 375.66			Driver Version: 375.66			
GPU	Name	Persistence-M	Bus-Id	Disp.A	Volatile Uncorr.	ECC
Fan	Temp	Pwr:Usage/Cap	Memory-Usage	GPU-Util	Compute M.	
0	Tesla K80	On	0000:00:17.0	Off	0	
N/A	59C	P0	97W / 149W	998MiB / 11439MiB	84%	Default
1	Tesla K80	On	0000:00:18.0	Off	0	
N/A	48C	P0	106W / 149W	826MiB / 11439MiB	80%	Default
2	Tesla K80	On	0000:00:19.0	Off	0	
N/A	60C	P0	92W / 149W	755MiB / 11439MiB	83%	Default
3	Tesla K80	On	0000:00:1A.0	Off	0	
N/A	52C	P0	114W / 149W	577MiB / 11439MiB	81%	Default
4	Tesla K80	On	0000:00:1B.0	Off	0	
N/A	60C	P0	92W / 149W	555MiB / 11439MiB	76%	Default
5	Tesla K80	On	0000:00:1C.0	Off	0	
N/A	48C	P0	107W / 149W	546MiB / 11439MiB	83%	Default
6	Tesla K80	On	0000:00:1D.0	Off	0	
N/A	63C	P0	92W / 149W	547MiB / 11439MiB	56%	Default
7	Tesla K80	On	0000:00:1E.0	Off	0	
N/A	51C	P0	103W / 149W	536MiB / 11439MiB	63%	Default
Processes:						
GPU	PID	Type	Process name	GPU Memory Usage		
1	55844	C	python	994MiB		
2	55844	C	python	822MiB		
3	55844	C	python	751MiB		
4	55844	C	python	573MiB		
5	55844	C	python	551MiB		
6	55844	C	python	542MiB		
7	55844	C	python	543MiB		
				532MiB		



33:45 / 42:33



# Amazon SageMaker



- A **fully managed** service that enables **data scientists** and **developers** to quickly and easily **build production-ready** machine-learning models.
- Amazon SageMaker includes several built-in state of the art algorithms that are fine-tuned to run on distributed environments.  
<https://docs.aws.amazon.com/sagemaker/latest/dg/algos.html>
- Amongst built-in algorithms there is a Sequence2Sequence algorithm implemented using Sockeye.
- All you need to do it to tune model's hyperparameters and passing your data to the model.
- Another useful built-in algorithm is Amazon BlazingText and provides Word2Vec.

```
In [7]: %%bash  
# python3 create_vocab_proto.py -h
```

The cell below does the preprocessing. If you are using the complete dataset, the script might take around 10-15 min on an m4.xlarge notebook instance.  
Remove ".small" from the file names for training on full datasets.

```
In [8]: %%time  
%%bash  
python3 create_vocab_proto.py \  
    --train-source corpus.tc.en.small \  
    --train-target corpus.tc.de.small \  
    --val-source validation/newstest2014.tc.en \  
    --val-target validation/newstest2014.tc.de  
  
INFO:__main__:Building vocabulary from dataset: corpus.tc.en.small and corpus.tc.de.small  
INFO:__main__:Final vocabulary: 10661 types (min frequency 1, top 50000 types)  
INFO:__main__:Final vocabulary: 15980 types (min frequency 1, top 50000 types)  
INFO:__main__:Source vocabulary size: 10661  
INFO:__main__:Vocabulary saved to "vocab.src.json"  
INFO:__main__:Target vocabulary size: 15980  
INFO:__main__:Vocabulary saved to "vocab.trg.json"  
INFO:__main__:Spawning 31 encoding worker(s) for encoding train datasets!  
INFO:__main__:Processed 10000 lines for encoding to protobuf. 0 lines were ignored as they didn't have  
any content in either the source or the target file!  
INFO:__main__:Completed writing the encoding queue!  
INFO:__main__:Encoding finished! Writing records to "train.rec"  
INFO:__main__:Processed input and saved to "train.rec"  
INFO:__main__:Spawning 31 encoding worker(s) for encoding validation datasets!  
INFO:__main__:Processed 3003 lines for encoding to protobuf. 0 lines were ignored as they didn't have  
any content in either the source or the target file!  
INFO:__main__:Completed writing the encoding queue!  
INFO:__main__:Encoding finished! Writing records to "val.rec"  
INFO:__main__:Processed input and saved to "val.rec"  
  
CPU times: user 0 ns, sys: 4 ms, total: 4 ms  
Wall time: 1.8 s
```



37:22 / 42:33



The script will output 4 files, namely:

- train.rec : Contains source and target sentences for training in protobuf format

```
INFO:__main__:Spawning 31 encoding worker(s) for encoding validation datasets!
INFO:__main__:Processed 3003 lines for encoding to protobuf. 0 lines were ignored as they didn't have
any content in either the source or the target file!
INFO:__main__:Completed writing the encoding queue!
INFO:__main__:Encoding finished! Writing records to "val.rec"
INFO:__main__:Processed input and saved to "val.rec"
```

```
CPU times: user 0 ns, sys: 4 ms, total: 4 ms
Wall time: 1.8 s
```

The script will output 4 files, namely:

- train.rec : Contains source and target sentences for training in protobuf format
- val.rec : Contains source and target sentences for validation in protobuf format
- vocab.src.json : Vocabulary mapping (string to int) for source language (English in this example)
- vocab.trg.json : Vocabulary mapping (string to int) for target language (German in this example)

Let's upload the pre-processed dataset and vocabularies to S3

```
In [11]: def upload_to_s3(bucket, prefix, channel, file):
    s3 = boto3.resource('s3')
    data = open(file, "rb")
    key = prefix + "/" + channel + '/' + file
    s3.Bucket(bucket).put_object(Key=key, Body=data)

upload_to_s3(bucket, prefix, 'train', 'train.rec')
upload_to_s3(bucket, prefix, 'validation', 'val.rec')
upload_to_s3(bucket, prefix, 'vocab', 'vocab.src.json')
upload_to_s3(bucket, prefix, 'vocab', 'vocab.trg.json')

In [12]: region_name = boto3.Session().region_name

In [13]: containers = {'us-west-2': '433757028032.dkr.ecr.us-west-2.amazonaws.com/seq2seq:latest',
                    'us-east-2': '825641698319.dkr.ecr.us-east-2.amazonaws.com/seq2seq:latest',
                    'eu-west-1': '685385470294.dkr.ecr.eu-west-1.amazonaws.com/seq2seq:latest'}
                    container = containers[region_name]
                    print('Using SageMaker Seq2Seq container: {} {}'.format(container, region_name))
```

37:33 / 42:33



Using SageMaker Seq2Seq container: 685385470294.dkr.ecr.eu-west-1.amazonaws.com/seq2seq:latest (eu-west-1)

- train.rec : Contains source and target sentences for training in protobuf format
- val.rec : Contains source and target sentences for validation in protobuf format
- vocab.src.json : Vocabulary mapping (string to int) for source language (English in this example)
- vocab.trg.json : Vocabulary mapping (string to int) for target language (German in this example)

Let's upload the pre-processed dataset and vocabularies to S3

```
In [11]: def upload_to_s3(bucket, prefix, channel, file):
    s3 = boto3.resource('s3')
    data = open(file, "rb")
    key = prefix + "/" + channel + '/' + file
    s3.Bucket(bucket).put_object(Key=key, Body=data)

upload_to_s3(bucket, prefix, 'train', 'train.rec')
upload_to_s3(bucket, prefix, 'validation', 'val.rec')
upload_to_s3(bucket, prefix, 'vocab', 'vocab.src.json')
upload_to_s3(bucket, prefix, 'vocab', 'vocab.trg.json')

In [12]: region_name = boto3.Session().region_name

In [13]: containers = {'us-west-2': '433757028032.dkr.ecr.us-west-2.amazonaws.com/seq2seq:latest',
                  'us-east-1': '811284229777.dkr.ecr.us-east-1.amazonaws.com/seq2seq:latest',
                  'us-east-2': '825641698319.dkr.ecr.us-east-2.amazonaws.com/seq2seq:latest',
                  'eu-west-1': '685385470294.dkr.ecr.eu-west-1.amazonaws.com/seq2seq:latest'}
    container = containers[region_name]
    print('Using SageMaker Seq2Seq container: {} ({})'.format(container, region_name))

Using SageMaker Seq2Seq container: 685385470294.dkr.ecr.eu-west-1.amazonaws.com/seq2seq:latest (eu-west-1)
```

## Training the Machine Translation model

```
In [14]: job_name = 'DEMO-seq2seq-en-de-' + strftime("%Y-%m-%d-%H", gmtime())
print("Training job", job_name)

create_training_params = \
{
    "AlgorithmSpecification": {
        "TrainingImage": container,
```

```
In [14]: job_name = 'DEMO-seq2seq-en-de-' + strftime("%Y-%m-%d-%H", gmtime())
print("Training job", job_name)

create_training_params = \
{
    "AlgorithmSpecification": {
        "TrainingImage": container,
        "TrainingInputMode": "File"
    },
    "RoleArn": role,
    "OutputDataConfig": {
        "S3OutputPath": "s3://{}{}".format(bucket, prefix)
    },
    "ResourceConfig": {
        # Seq2Seq does not support multiple machines. Currently, it only supports single machine, multiple GPUs
        "InstanceCount": 1,
        "InstanceType": "ml.p3.8xlarge", # We suggest one of ["ml.p2.16xlarge", "ml.p2.8xlarge", "ml.p2.xlarge"]
        "VolumeSizeInGB": 50
    },
    "TrainingJobName": job_name,
    "HyperParameters": {
        # Please refer to the documentation for complete list of parameters
        "max_seq_len_source": "60",
        "max_seq_len_target": "60",
        "optimized_metric": "bleu",
        "batch_size": "64", # Please use a larger batch size (256 or 512) if using ml.p2.8xlarge or ml.p2.16xlarge
        "checkpoint_frequency_num_batches": "1000",
        "rnn_num_hidden": "512",
        "num_layers_encoder": "1",
        "num_layers_decoder": "1",
        "num_embed_source": "512",
        "num_embed_target": "512",
        "checkpoint_threshold": "3",
        "max_num_batches": "2100"
        # Training will stop after 2100 iterations/batches.
        # This is just for demo purposes. Remove the above parameter if you want a better model.
    },
    "StoppingCondition": {
        "MaxRuntimeInSeconds": 48 * 3600
    },
    "InputDataConfig": [
        {
            "DataSource": {
                "S3DataSource": {
                    "S3Uri": "s3://{}{}".format(bucket, prefix),
                    "ContentType": "Text"
                }
            },
            "RecordWrapperType": "Line"
        }
    ]
}
```



37:54 / 42:33



```
primary_container = {
    'Image': container,
    'ModelDataUrl': model_data
}

create_model_response = sage.create_model(
    ModelName = model_name,
    ExecutionRoleArn = role,
    PrimaryContainer = primary_container)

print(create_model_response['ModelArn'])

DEMO-seq2seq-en-de-2018-05-03-23
s3://cyrusmv-sagemaker-demos/sagemaker/DEMO-seq2seq/DEMO-seq2seq-en-de-2018-05-03-23/output/model.tar.gz
arn:aws:sagemaker:eu-west-1:XXXXXXXXXXXXX:model/demo-seq2seq-en-de-2018-05-03-23
CPU times: user 16 ms, sys: 0 ns, total: 16 ms
Wall time: 496 ms
```

## Create endpoint configuration

Use the model to create an endpoint configuration. The endpoint configuration also contains information about the type and number of EC2 instances to use when hosting the model.

Since SageMaker Seq2Seq is based on Neural Nets, we could use an ml.p2.xlarge (GPU) instance, but for this example we will use a free tier eligible ml.m4.xlarge.

```
In [19]: from time import gmtime, strftime

endpoint_config_name = 'DEMO-Seq2SeqEndpointConfig-' + strftime("%Y-%m-%d-%H-%M-%S", gmtime())
print(endpoint_config_name)
create_endpoint_config_response = sage.create_endpoint_config(
    EndpointConfigName = endpoint_config_name,
    ProductionVariants=[{
        'InitialInstanceCount':1,
        'ModelName':model_name,
        'VariantName':'AllTraffic'}])

print('Endpoint Config Arn: ' + create_endpoint_config_response['EndpointConfigArn'])
```

## Create endpoint configuration

Use the model to create an endpoint configuration. The endpoint configuration also contains information about the type and number of EC2 instances to use when hosting the model.

Since SageMaker Seq2Seq is based on Neural Nets, we could use an ml.p2.xlarge (GPU) instance, but for this example we will use a free tier eligible ml.m4.xlarge.

```
In [19]: from time import gmtime, strftime

endpoint_config_name = 'DEMO-Seq2SeqEndpointConfig-' + strftime("%Y-%m-%d-%H-%M-%S", gmtime())
print(endpoint_config_name)
create_endpoint_config_response = sage.create_endpoint_config(
    EndpointConfigName = endpoint_config_name,
    ProductionVariants=[{
        'InstanceType': 'ml.m4.xlarge',
        'InitialInstanceCount': 1,
        'ModelName': model_name,
        'VariantName': 'AllTraffic'})]

print("Endpoint Config Arn: " + create_endpoint_config_response['EndpointConfigArn'])

DEMO-Seq2SeqEndpointConfig-2018-05-03-46-48
Endpoint Config Arn: arn:aws:sagemaker:eu-west-1: :endpoint-config/demo-seq2seqendpointconfig-2018-05-03-2
3-46-48
```

## Create endpoint

Lastly, we create the endpoint that serves up model, through specifying the name and configuration defined above. The end result is an endpoint that can be validated and incorporated into production applications. This takes 10-15 minutes to complete.

```
In [20]: #%%time
import time

endpoint_name = 'DEMO-Seq2SeqEndpoint-' + strftime("%Y-%m-%d-%H-%M-%S", gmtime())
print(endpoint_name)
create_endpoint_response = sage.create_endpoint(
    EndpointName=endpoint_name,
```



39:22 / 42:23



```
initialInstanceCount :1,
'ModelName':model_name,
'VariantName':'AllTraffic'}})

print("Endpoint Config Arn: " + create_endpoint_config_response['EndpointConfigArn'])

DEMO-Seq2SeqEndpointConfig-2018-05-03-23-46-48
Endpoint Config Arn: arn:aws:sagemaker:eu-west-1:.....:endpoint-config/demo-seq2seqendpointconfig-2018-05-03-2
3-46-48
```

## Create endpoint

Lastly, we create the endpoint that serves up model, through specifying the name and configuration defined above. The end result is an endpoint that can be validated and incorporated into production applications. This takes 10-15 minutes to complete.

```
In [20]: %%time
import time

endpoint_name = 'DEMO-Seq2SeqEndpoint-' + strftime("%Y-%m-%d-%H-%M-%S", gmtime())
print(endpoint_name)
create_endpoint_response = sage.create_endpoint(
    EndpointName=endpoint_name,
    EndpointConfigName=endpoint_config_name)
print(create_endpoint_response['EndpointArn'])

resp = sage.describe_endpoint(EndpointName=endpoint_name)
status = resp['EndpointStatus']
print("Status: " + status)

# wait until the status has changed
sage.get_waiter('endpoint_in_service').wait(EndpointName=endpoint_name)

# print the status of the endpoint
endpoint_response = sage.describe_endpoint(EndpointName=endpoint_name)
status = endpoint_response['EndpointStatus']
print('Endpoint creation ended with EndpointStatus = {}'.format(status))

if status != 'InService':
    raise Exception('Endpoint creation failed.')
```

## Create endpoint

Lastly, we create the endpoint that serves up model, through specifying the name and configuration defined above. The end result is an endpoint that can be validated and incorporated into production applications. This takes 10-15 minutes to complete.

```
In [20]: #time
import time

endpoint_name = 'DEMO-Seq2SeqEndpoint-' + strftime("%Y-%m-%d-%H-%M-%S", gmtime())
print(endpoint_name)
create_endpoint_response = sage.create_endpoint(
    EndpointName=endpoint_name,
    EndpointConfigName=endpoint_config_name)
print(create_endpoint_response['EndpointArn'])

resp = sage.describe_endpoint(EndpointName=endpoint_name)
status = resp['EndpointStatus']
print("Status: " + status)

# wait until the status has changed
sage.get_waiter('endpoint_in_service').wait(EndpointName=endpoint_name)

# print the status of the endpoint
endpoint_response = sage.describe_endpoint(EndpointName=endpoint_name)
status = endpoint_response['EndpointStatus']
print('Endpoint creation ended with EndpointStatus = {}'.format(status))

if status != 'InService':
    raise Exception('Endpoint creation failed.')

DEMO-Seq2SeqEndpoint-2018-05-03-23-46-52
arn:aws:sagemaker:eu-west-1:... :endpoint/demo-seq2segendpoint-2018-05-03-23-46-52
Status: Creating
Endpoint creation ended with EndpointStatus = InService
CPU times: user 56 ms, sys: 0 ns, total: 56 ms
Wall time: 6min 31s
```

If you see the message,

Endpoint creation ended with EndpointStatus = InService

then congratulations! You now have a functioning inference endpoint. You can confirm the endpoint configuration and status by navigating to the "Endpoints" tab in the AWS SageMaker console.

We will finally create a runtime object from which we can invoke the endpoint.

```
In [21]: runtime = boto3.client(service_name='runtime.sagemaker')
```

## Perform Inference

### Using JSON format for inference (Suggested for a single or small number of data instances)

Note that you don't have to convert string to text using the vocabulary mapping for inference using JSON mode

```
In [22]: sentences = ["you are so good !",
                    "can you drive a car ?",
                    "i want to watch a movie ."
                   ]

payload = {"instances" : []}
for sent in sentences:
    payload["instances"].append({"data" : sent})

response = runtime.invoke_endpoint(EndpointName=endpoint_name,
                                   ContentType='application/json',
                                   Body=json.dumps(payload))

response = response["Body"].read().decode("utf-8")
response = json.loads(response)
print(response)

{'predictions': [{'target': 'COR / 08 / 08 / 149'}, {'target': 'wir werden .'}, {'target': 'Präsident Maystadt sollte\nn sich : '}])
```

## Perform Inference

Using JSON format for inference (Suggested for a single or small number of data instances)

Note that you don't have to convert string to text using the vocabulary mapping for inference using JSON mode

```
In [22]: sentences = ["you are so good !",
                    "can you drive a car ?",
                    "i want to watch a movie ."
                   ]

payload = {"instances" : []}
for sent in sentences:
    payload["instances"].append({"data" : sent})

response = runtime.invoke_endpoint(EndpointName=endpoint_name,
                                   ContentType='application/json',
                                   Body=json.dumps(payload))

response = response["Body"].read().decode("utf-8")
response = json.loads(response)
print(response)

{'predictions': [{'target': 'COR / 08 / 08 / 149'}, {'target': 'wir werden .'}, {'target': 'Präsident Maystadt sollte\nn sich : '}]}
```

## Retrieving the Attention Matrix

Passing "attention\_matrix":"true" in configuration of the data instance will return the attention matrix.

```
In [23]: sentence = 'can you drive a car ?'

payload = {"instances" : [
                           {"data" : sentence}
```

```
{'predictions': [{target': 'COR / 08 / 08 / 149'}, {'target': 'wir werden .'}, {'target': "Präsident Maystadt sollte\nn sich : "}]}
```

## Retrieving the Attention Matrix

Passing `"attention_matrix": "true"` in configuration of the data instance will return the attention matrix.

```
In [23]: sentence = 'can you drive a car ?'

payload = {"instances" : [
    {
        "data" : sentence,
        "configuration" : {"attention_matrix": "true"}
    }
]}

response = runtime.invoke_endpoint(EndpointName=endpoint_name,
                                   ContentType='application/json',
                                   Body=json.dumps(payload))

response = response["Body"].read().decode("utf-8")
response = json.loads(response)['predictions'][0]

source = sentence
target = response["target"]
attention_matrix = np.array(response["matrix"])

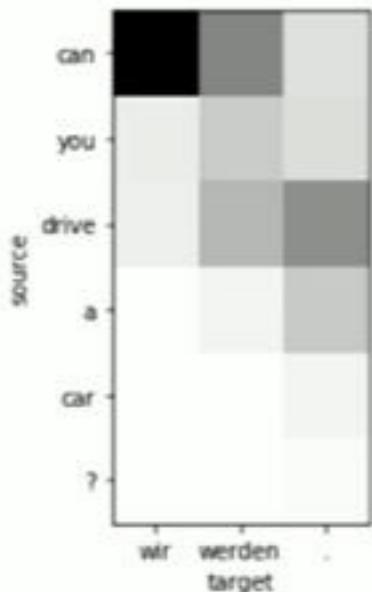
print("Source: %s \nTarget: %s" % (source, target))
```

Source: can you drive a car ?  
Target: wir werden .

```
In [24]: # Define a function for plotting the attention matrix
def plot_matrix(attention_matrix, target, source):
    source_tokens = source.split()
    target_tokens = target.split()
    assert attention_matrix.shape[0] == len(target_tokens)
    plt.imshow(attention_matrix.transpose(), interpolation="nearest", cmap="Greys")
    plt.xlabel("target")
    plt.ylabel("source")
```

```
plt.xlabel("target")
plt.ylabel("source")
plt.gca().set_xticks([i for i in range(0, len(target_tokens))])
plt.gca().set_yticks([i for i in range(0, len(source_tokens))])
plt.gca().set_xticklabels(target_tokens)
plt.gca().set_yticklabels(source_tokens)
plt.tight_layout()
```

```
In [25]: plot_matrix(attention_matrix, target, source)
```



### Using Protobuf format for inference (Suggested for efficient bulk inference)

Reading the vocabulary mappings as this mode of inference accepts list of integers and returns list of integers.

```
In [26]: import io
import tempfile
from record_pb2 import Record
from create_vocab_proto import vocab_from_json, reverse_vocab, write_recordio, list_to_record_bytes, read_next

source = vocab_from_json("vocab.src.json")
target = vocab_from_json("vocab.trg.json")
```

[What Is Amazon SageMaker?](#)[How It Works](#)[Getting Started](#)[Using Built-in Algorithms](#)[Common Information](#)[Linear Learner](#)[Factorization Machines](#)[XGBoost Algorithm](#)[Image Classification Algorithm](#)[Sequence to Sequence \(seq2seq\)](#)[How It Works](#)[Hyperparameters](#)[K-Means Algorithm](#)[Principal Component Analysis \(PCA\)](#)[Latent Dirichlet Allocation \(LDA\)](#)[Neural Topic Model \(NTM\)](#)[DeepAR Forecasting](#)[BlazingText](#)[Random Cut Forest](#)[Using Your Own Algorithms](#)[Automatically Scaling Amazon SageMaker Models](#)[Using TensorFlow](#)[Using Apache MXNet](#)

## Sequence-to-sequence hyperparameters

Parameter Name	Description
max_seq_len_source	Maximum length for the source sequence length. Sequences longer than this length are truncated to this length. Valid values: positive integer Default value: 100
max_seq_len_target	Maximum length for the target sequence length. Sequences longer than this length are truncated to this length. Valid values: positive integer Default value: 100
encoder_type	Encoder type. The <i>rnn</i> architecture is based on attention mechanism by Bahdanau et al. and <i>cnn</i> architecture is based on Gehring et al. Valid values: String. Either <i>rnn</i> or <i>cnn</i> . Default value: <i>rnn</i>
decoder_type	Decoder type. Valid values: String. Either <i>rnn</i> or <i>cnn</i> . Default value: <i>rnn</i>
num_layers_encoder	Number of layers for Encoder <i>rnn</i> or <i>cnn</i> . Valid values: positive integer Default value: 1
num_layers_decoder	Number of layers for Decoder <i>rnn</i> or <i>cnn</i> . Valid values: positive integer Default value: 1
rnn_num_hidden	The number of <i>rnn</i> hidden units for encoder and decoder. This must be a multiple of 2 because the algorithm uses bi-directional LSTM by default. Valid values: positive integer

[What Is Amazon SageMaker?](#)[How It Works](#)[Getting Started](#)[Using Built-in Algorithms](#)[Common Information](#)[Linear Learner](#)[Factorization Machines](#)[XGBoost Algorithm](#)[Image Classification Algorithm](#)[Sequence to Sequence \(seq2seq\)](#)[How It Works](#)[Hyperparameters](#)[K-Means Algorithm](#)[Principal Component Analysis \(PCA\)](#)[Latent Dirichlet Allocation \(LDA\)](#)[Neural Topic Model \(NTM\)](#)[DeepAR Forecasting](#)[BlazingText](#)[Random Cut Forest](#)[Using Your Own Algorithms](#)[Automatically Scaling Amazon SageMaker Models](#)[TensorFlow](#)[Apache MXNet](#)

42:11 / 42:33



English

Sign In to the Console

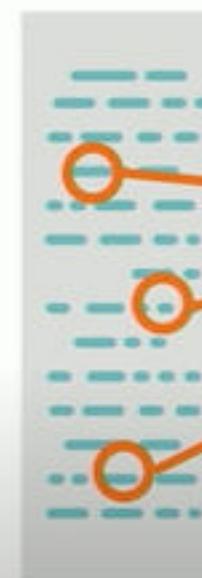
		Valid values: positive integer Default value: 5
	cnn_num_hidden	Number of <i>cnn</i> hidden units for encoder and decoder. Valid values: positive integer Default value: 512
	cnn_activation_type	The <i>cnn</i> activation type to be used. Valid values: String. One of <i>glu</i> , <i>relu</i> , <i>softrelu</i> , <i>sigmoid</i> , or <i>tanh</i> . Default value: <i>glu</i>
	cnn_hidden_dropout	Dropout probability for dropout between convolutional layers. Valid values: Float. Range in [0,1]. Default value: 0
	num_embed_source	Embedding size for source tokens. Valid values: positive integer Default value: 512
	num_embed_target	Embedding size for target tokens. Valid values: positive integer Default value: 512
	embed_dropout_source	Dropout probability for source side embeddings. Valid values: Float. Range in [0,1]. Default value: 0
	embed_dropout_target	Dropout probability for target side embeddings. Valid values: Float. Range in [0,1]. Default value: 0
	rnn_attention_type	Attention model for encoders. <i>mlp</i> refers to concat and bilinear refers to general from the Luong et al. paper.

- Neubig, G. (2017). *Neural Machine Translation and Sequence-to-Sequence Models: A Tutorial*. Carnegie Mellon University: Language Technologies Institute.
- Cliffton, A., Denkowski, M., Domhan, T., Hieber, F., Post, M., Sokolov, A., & Vilar, D. (2017). Sockeye Documentation. <http://sockeye.readthedocs.io/en/latest/index.html>
- Abadi, M., Agarwal, A., Barham, P., Brevdo, E., Chen, Z., Citro, C.,...Zheng, X. (2015) Vector Representations of Words. *Tutorials*.
- Abadi, M., Agarwal, A., Barham, P., Brevdo, E., Chen, Z., Citro, C.,...Zheng, X. (2015) Vector Representations of Words. *Github*.
- Weisstein, E. (2018). Wolfram MathWorld. *Affine Transformation*.
- Klein, G., Kim, Y., Deng, Y., Senellart, J. & Rush, A. (2018). OpenNMT: Open-Source Toolkit for Neural Machine Translation. *ArXiv e-prints*.
- Baumel, T. (2018). Sequence to Sequence Attention Models in DyNet. *Github*.

# What problems exist that Comprehend Medical can solve?



## Electronic Health Records (EHR)

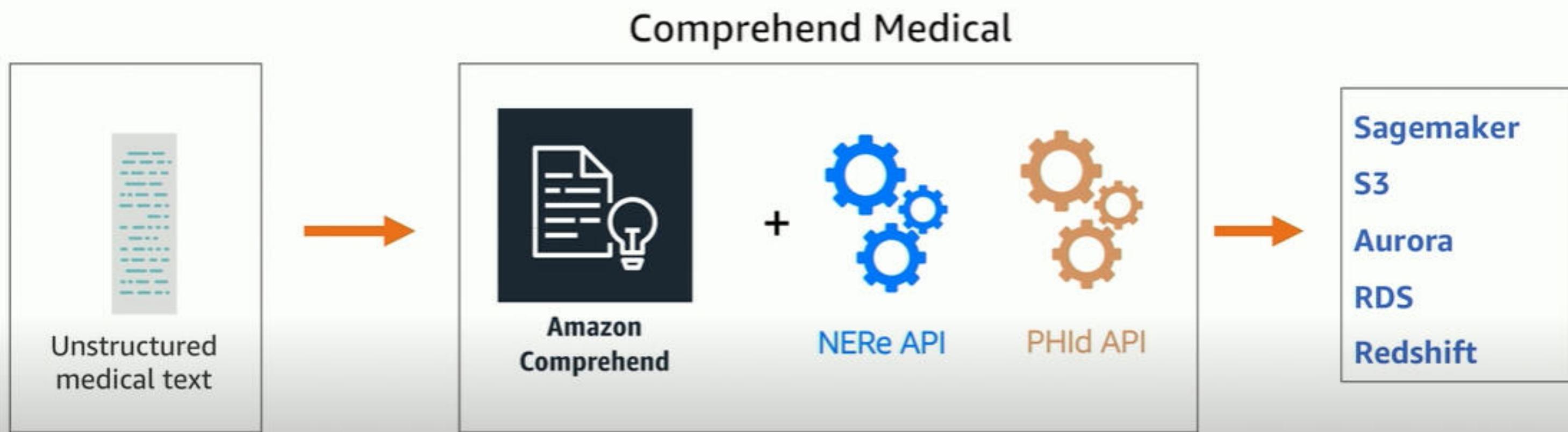


Extracting valuable patient-care information

Manual and labor intensive process



# How does Comprehend Medical Work?





Services

Resource Groups

Amazon Comprehend

Real-time analysis

Job management

## Customization

Custom classification

## Comprehend Medical

Real-time analysis

Beta

Comprehend Medical &gt; Real-time analysis

## Real-time analysis

Beta

See how Comprehend Medical recognizes entities related to the healthcare domain. To analyze your text, type or paste it in the text box.

## Input text

GENERAL : Well-developed , well-nourished male in no acute distress .  
HEENT : Oropharynx clear . Sclerae are anicteric .  
CARDIOVASCULAR : Regular rate and rhythm . No murmurs , rubs , or gallops .  
CHEST : Lungs CTA bilaterally .  
ABDOMEN : Benign .  
EXTREMITIES : No clubbing , cyanosis , or edema .

1309 of 20000 characters used.

[Clear text](#)[Analyze](#)

## Insights

## Analyzed text

## LABAHN HOSPITAL NOTE

• Address (LABAHN HOSPITAL)

## HISTORY OF PRESENT ILLNESS

Mr. Rickie Cresta is a 26 - year - old male with refractory ITP.  
• Name (Rickie Cresta) • Age (26)  
• Dx name (refractory ITP)

## INTERVAL HISTORY

Mr. Burl Bickler was last seen on July 11th 2016 , with Dr. Menzel Eplett . Mr. Rubin Corrado has  
• Name (Burl Bickler) • Date (July 11th 2016)  
• Name (Menzel Eplett) • Name (Rubin Corrado)

since started on romiplostim on May 13th 2007 . Despite the gradual dose increases in romiplostim ,  
• Generic name (romiplostim) • Date (May 13th 2007)  
• Generic name (romiplostim)

03:26 / 07:26

his platelet count has varied significantly , ranging anywhere from 5 to almost 60,000 . Recently Mr.



Test value



Services

Resource Groups



## Amazon Comprehend



Real-time analysis

Job management

## Customization

Custom classification

## Comprehend Medical

Beta

Real-time analysis

1309 of 20000 characters used.

Clear text

Analyze



## Insights

## Analyzed text

## PHYSICAL EXAMINATION

• Test name (PHYSICAL EXAMINATION)

**VITAL SIGNS :** Temperature is 36, pulse of 80, blood pressure 112 / 60.

• Test name (VITAL SIGNS) • Test name (Temperature) • Test value (36) • Test name (pulse) • Test value (80) • Test name (blood pressure) • Test value (112 / 60)

**GENERAL :** Well-developed, well-nourished male in no acute distress.

• System organ site (GENERAL) • Dx name (Well-developed) • Dx name (well-nourished) • Acuity (acute) • Dx name (distress)

**HEENT :** Oropharynx clear, Sclerae are anicteric.

• System organ site (HEENT) • System organ site (Oropharynx) • Dx name (Oropharynx clear) • System organ site (Sclerae) • Dx name (Sclerae are anicteric)

**CARDIOVASCULAR :** Regular rate and rhythm . No murmurs, rubs, or gallops.

• System organ site (CARDIOVASCULAR) • Dx name (Regular rate and rhythm) • Dx name (murmurs) • Dx name (rubs) • Dx name (gallops)

**CHEST :** Lungs CTA bilaterally.

• System organ site (CHEST) • System organ site (Lungs) • Direction (bilaterally) • Dx name (Lungs CTA)

We launch

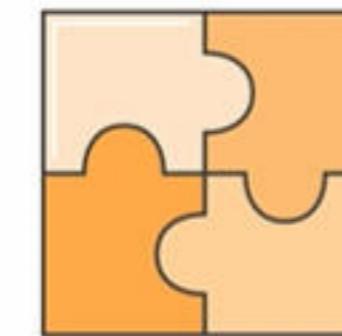
# Features of Amazon Translate



Secure  
Communication



Authentication and  
Access Control



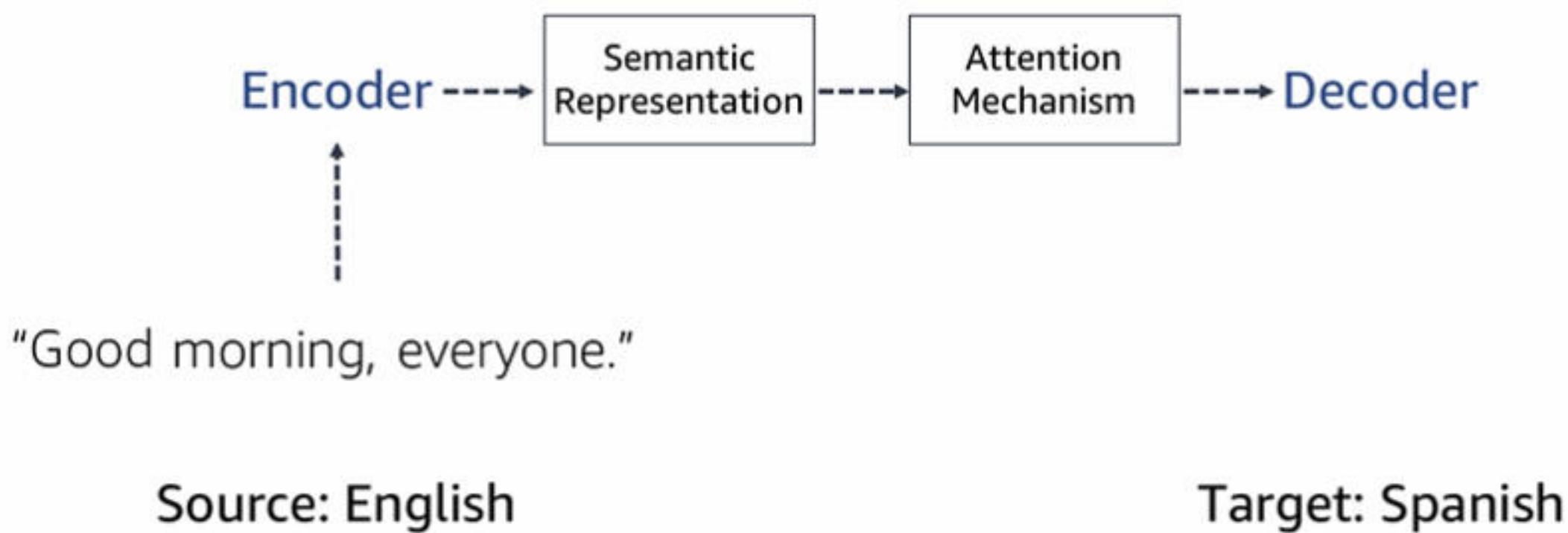
Integration with  
AWS Services

# Applications for Any Language

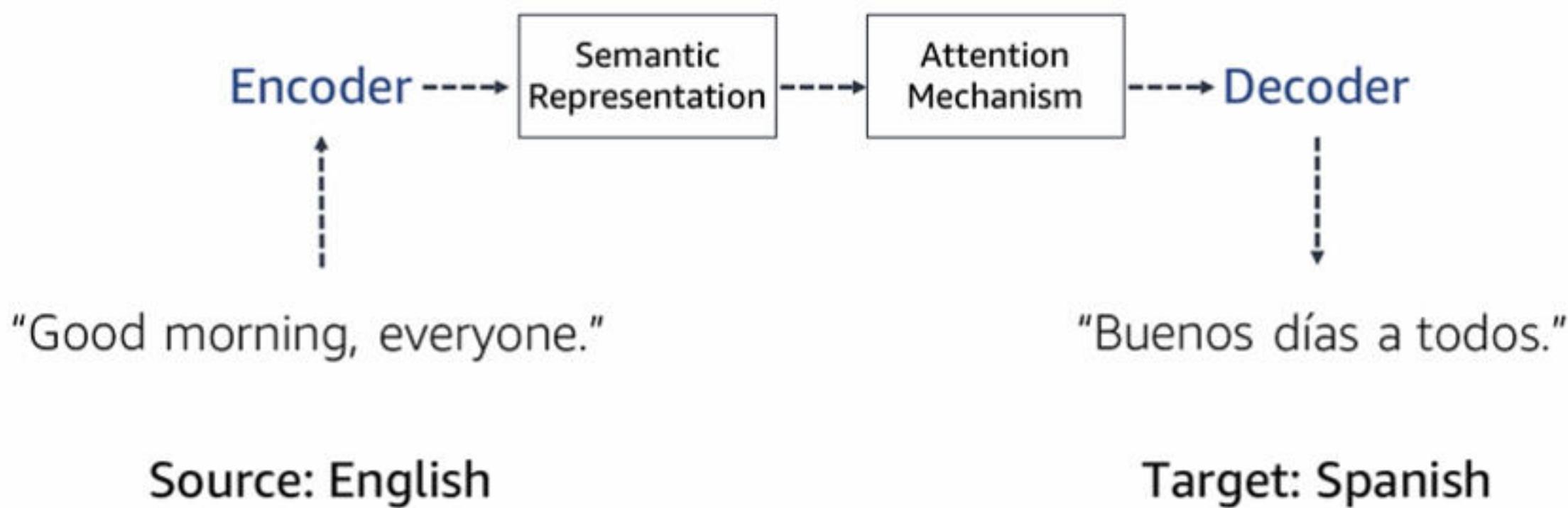
**aws** training and certification



# How It Works



# How It Works



# Getting Started



**1**

Set up  
AWS account

**2**

Create  
IAM user

**3**

Connect: Console,  
AWS CLI, AWS SDKs

# Demo: Connecting from the AWS Management Console

# Connecting from the AWS CLI



**Short text:** Input directly into command line

```
aws translate translate-text \
    --endpoint-url endpoint \
    --region region \
    --source-language-code "en" \
    --target-language-code "es" \
    --text "Hello, World"
```

**Long Text:** Use JSON file

```
aws translate translate-text \
    --endpoint-url endpoint \
    --region region \
    --cli-input-json file://translate.json > translated.json
```

# Using AWS SDKs



```
import boto3

translate = boto3.client(service_name='translate', region_name='region',
                         endpoint_url='endpoint', use_ssl=True)

result = translate.translate_text(Text="Hello, World",
                                   SourceLanguageCode="en", TargetLanguageCode="de")
print('TranslatedText: ' + result.get('TranslatedText'))
print('SourceLanguageCode: ' + result.get('SourceLanguageCode'))
print('TargetLanguageCode: ' + result.get('TargetLanguageCode'))
```



## Amazon Translate



Translating web-authored content  
in real time and on demand

Batch translating pre-existing  
content for analysis and insights

[Back To Products](#)

## 1966 SHELBY COBRA 427 S/C

Classic Cars

This model of the 1966 Shelby Cobra 427 S/C includes many authentic details and operating parts. The new car was designed in cooperation with Ford in Detroit. A new chassis was built using 4 in (102 mm) main chassis tubes (up from 3 in (76 mm)) and coil spring suspension all around. The new car also had wide fenders and a larger radiator opening. It was powered by the "side oiler" Ford 427 engine (7.0 L) rated at 425 bhp (317 kW), which provided a top speed of 164 mph (262 km/h) in the standard model and 485 bhp (362 kW) with a top speed of 185 mph (298 km/h) in the competition model.

CURRENT PRICE: \$74000

[ADD TO CART](#)

### REVIEWS

30/11/2017 17:30:03 Carlos Escapa wrote:

A mi tambien me gusta. Tendria que ir a la peluqueria todos los dias.

[Translate](#)



01/03/2018 06:55:32 Diego Natalli wrote:

I love this car a lot. I believe that James Bond had one. I wish I had one but I am not as rich as Bill Gates!!!

[Translate](#)



23/03/2018 12:00:19 Giuseppe Porcelli wrote:

When I bought it, I didn't like its size. So I hate it!

[Translate](#)



23/03/2018 12:01:16 Kris Skrinak wrote:

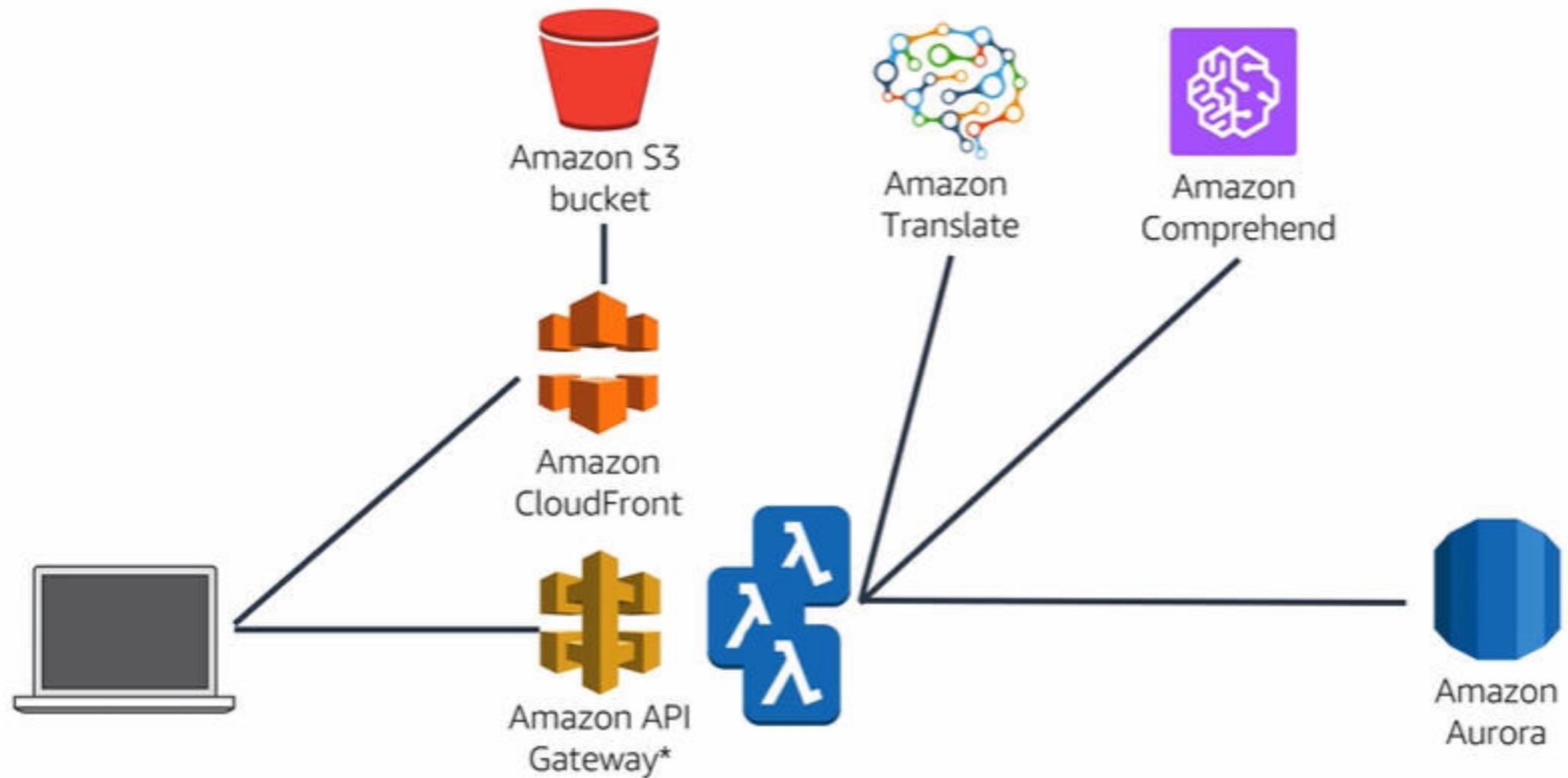
Slower than my Ducati, faster than my one wheel. Love it!!

[Translate](#)

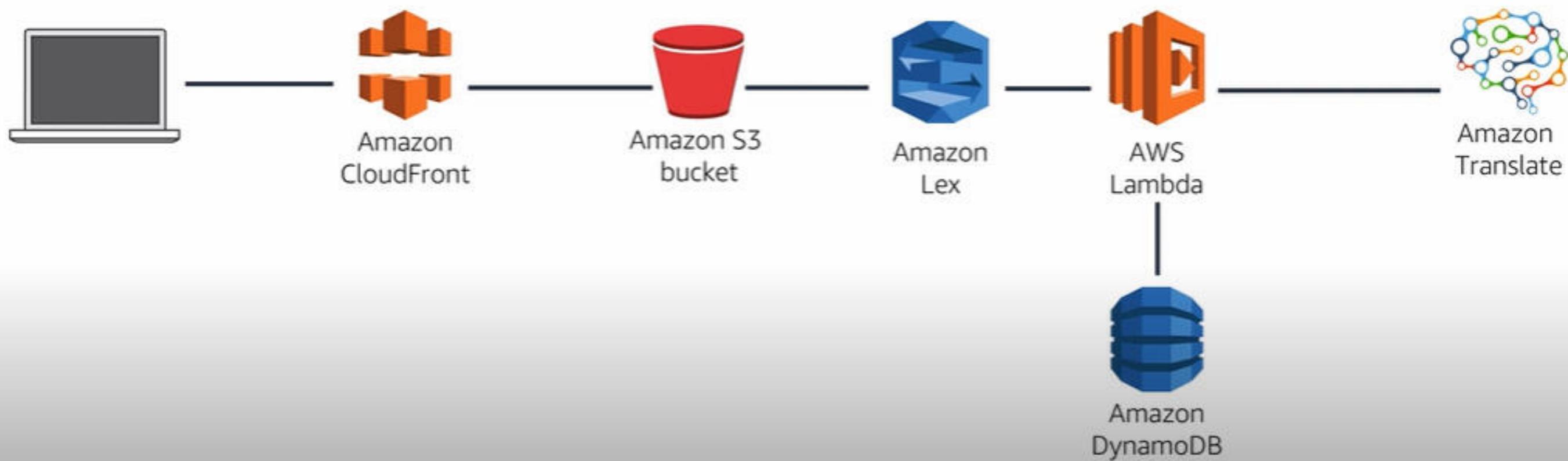


Please login to add or update your review

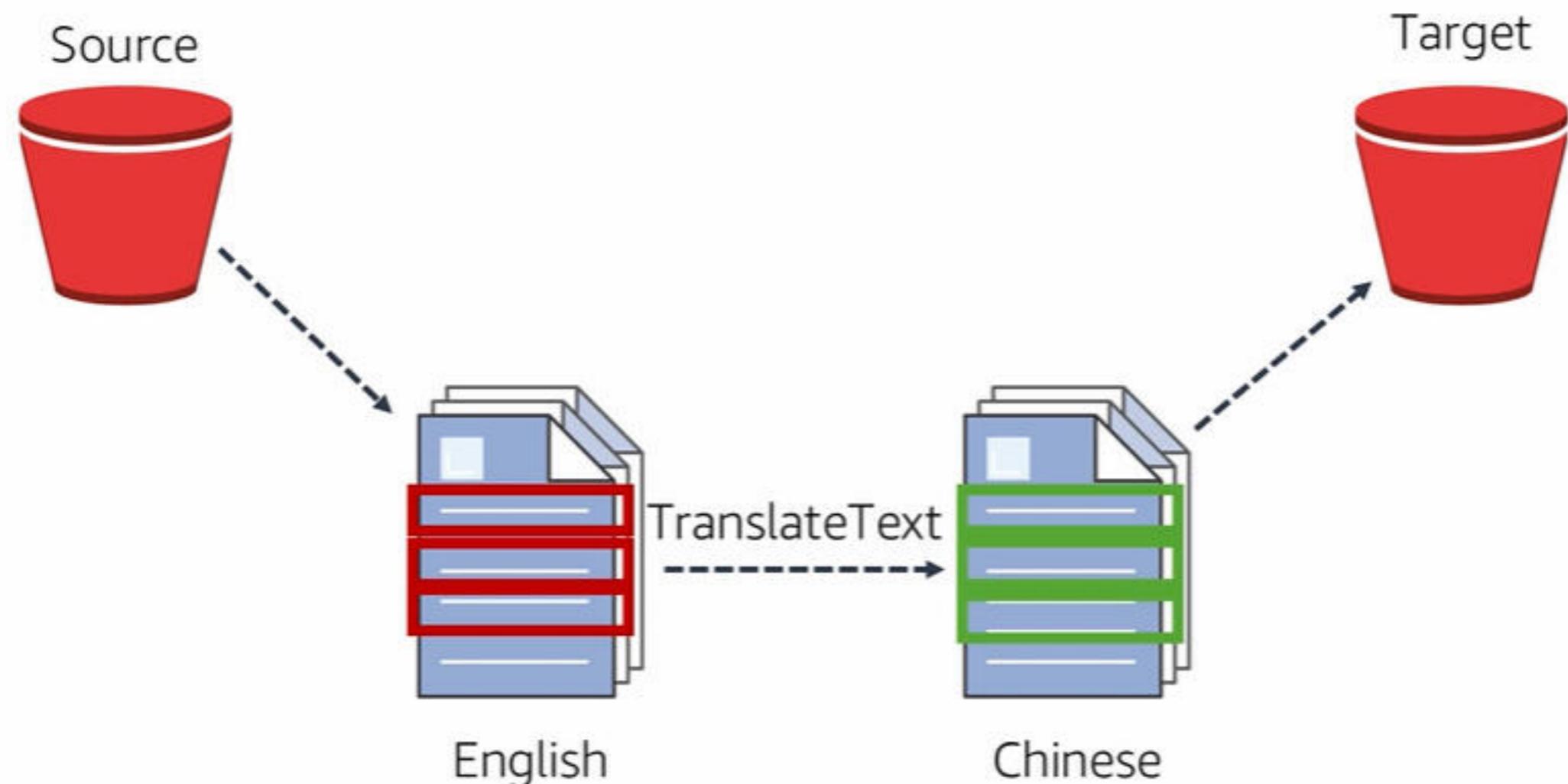
# Example: On-Demand Translation



# Example: Chatbot Translation



# Example: Batch Translations



- Neural machine translation engine
- Real-time, on-demand, and batch translations
- Integrates with AWS services
- Connect via the AWS Management Console, AWS CLI, and AWS SDKs



# Thanks for watching!

© 2018 Amazon Web Services, Inc. or its affiliates. All rights reserved. This work may not be reproduced or redistributed, in whole or in part, without prior written permission from Amazon Web Services, Inc. Commercial copying, lending, or selling is prohibited. Corrections or feedback on the course, please email us at: [aws-course-feedback@amazon.com](mailto:aws-course-feedback@amazon.com). For all other questions, contact us at: <https://aws.amazon.com/contact-us/aws-training/>. All trademarks are the property of their owners.

# Why Amazon Transcribe?



Readable, helpful transcriptions

Four score and seven years ago our fathers brought forth on this continent a new nation conceived in liberty and dedicated to the proposition that all men are created equal. Now we are engaged in a great civil war testing whether that nation or any nation so conceived and so dedicated can long endure.

```
[{"start_time": "0.150", "end_time": "0.580", "alternatives":  
[ {"confidence": "0.9984", "content": "Four"} ] ,
```

# Why Amazon Transcribe?

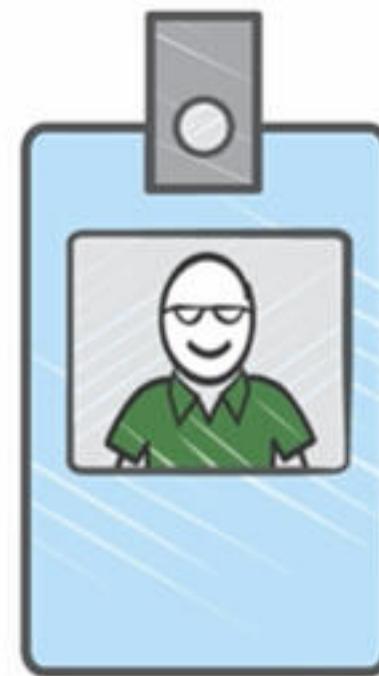


## Source-specific accuracy

- 큐 Custom vocabulary
- 큐 Language selection
- 큐 Multiple speaker recognition



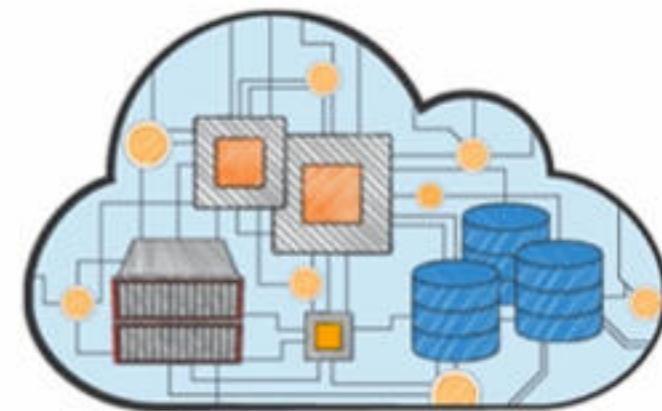
# Why Amazon Transcribe?



API Interface

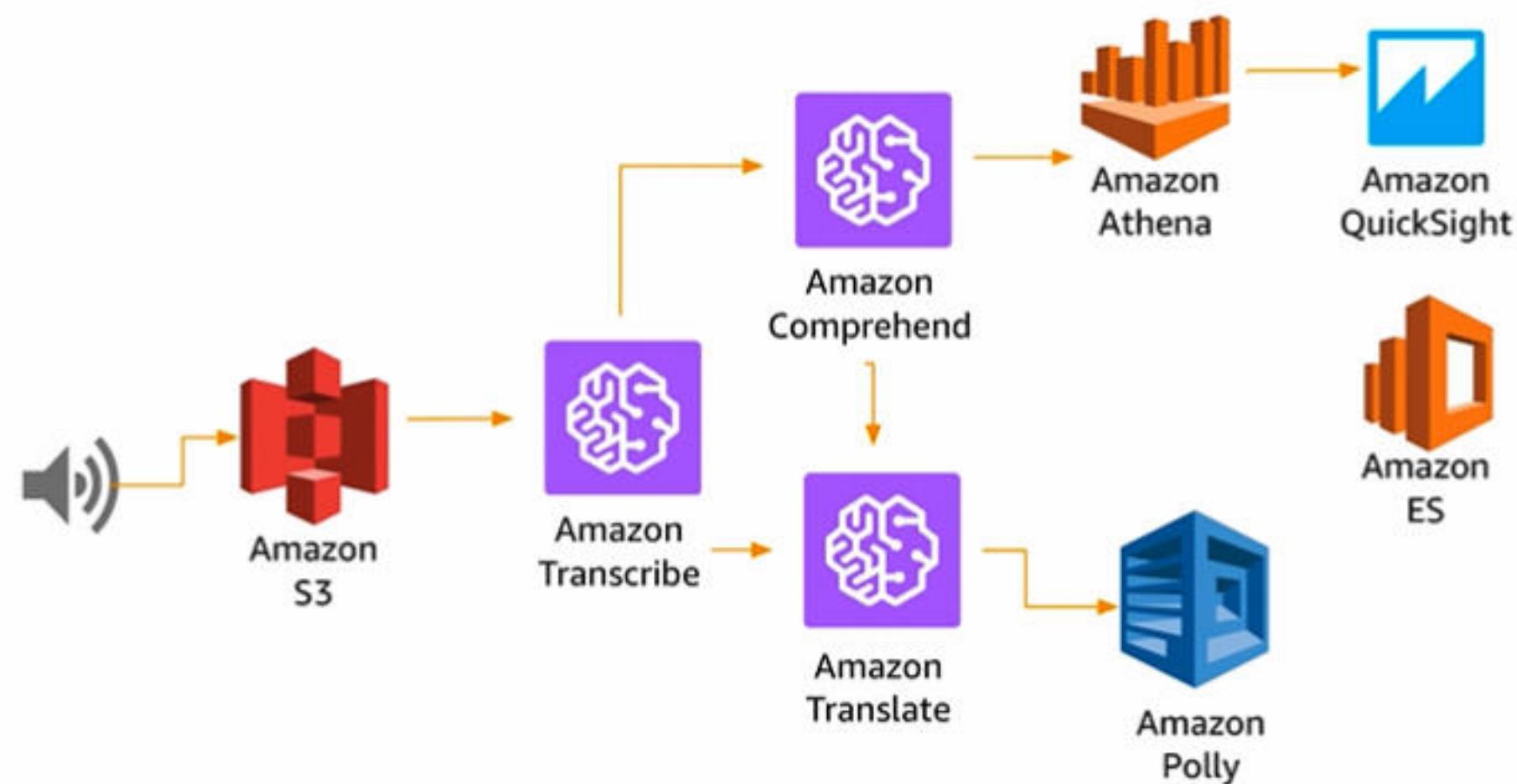


Authentication  
and Access  
Control

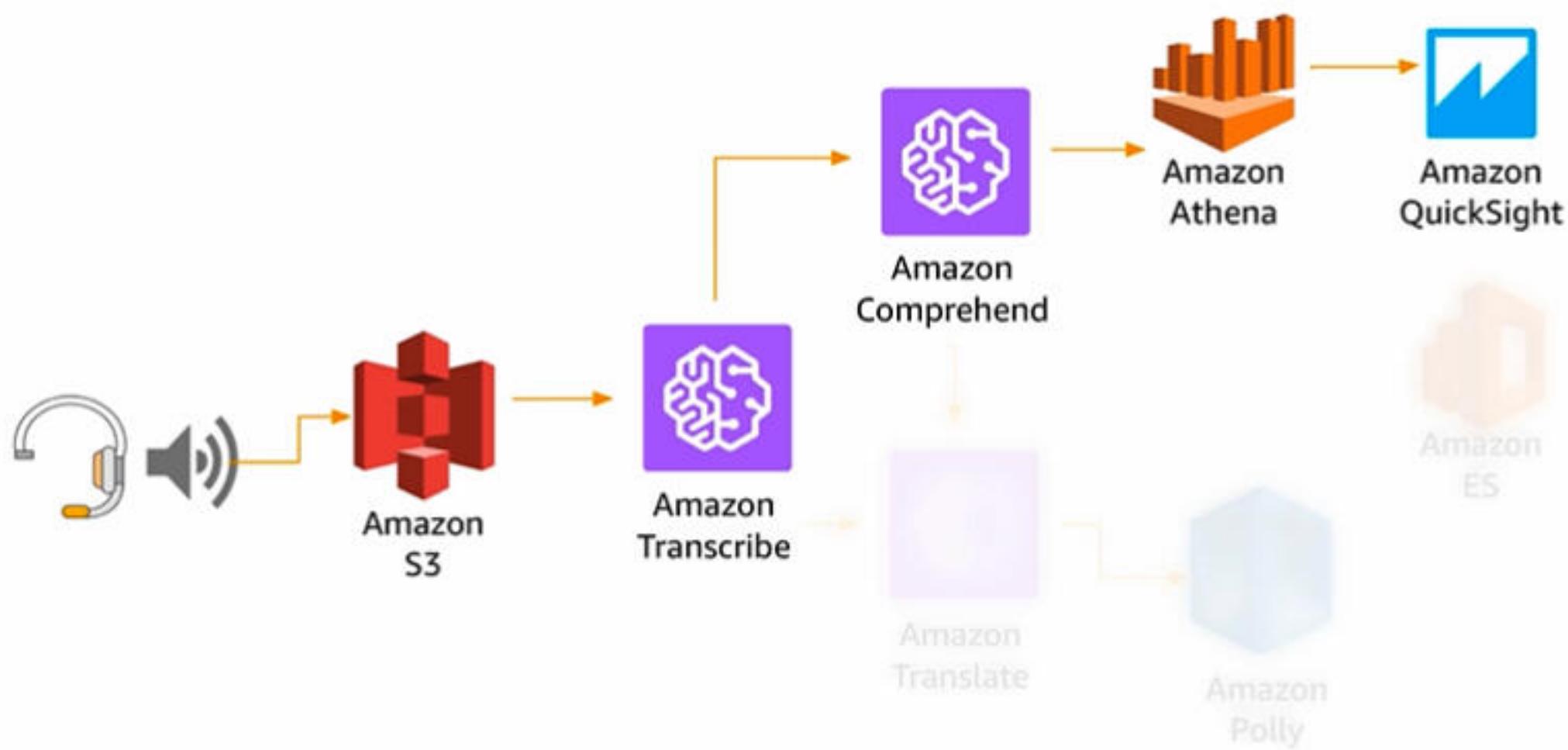


Integration with  
AWS Services

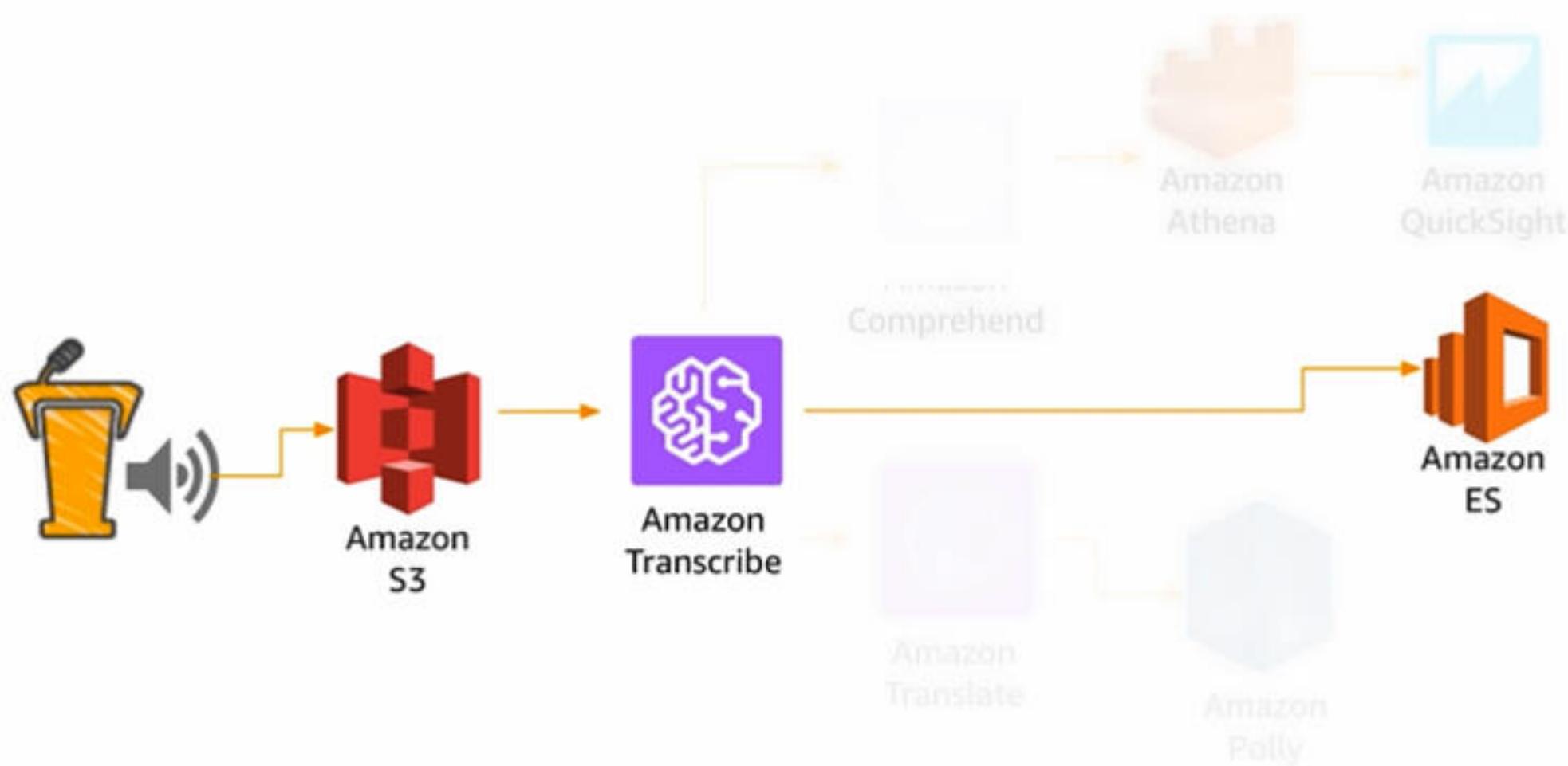
# Integrating with AWS Services



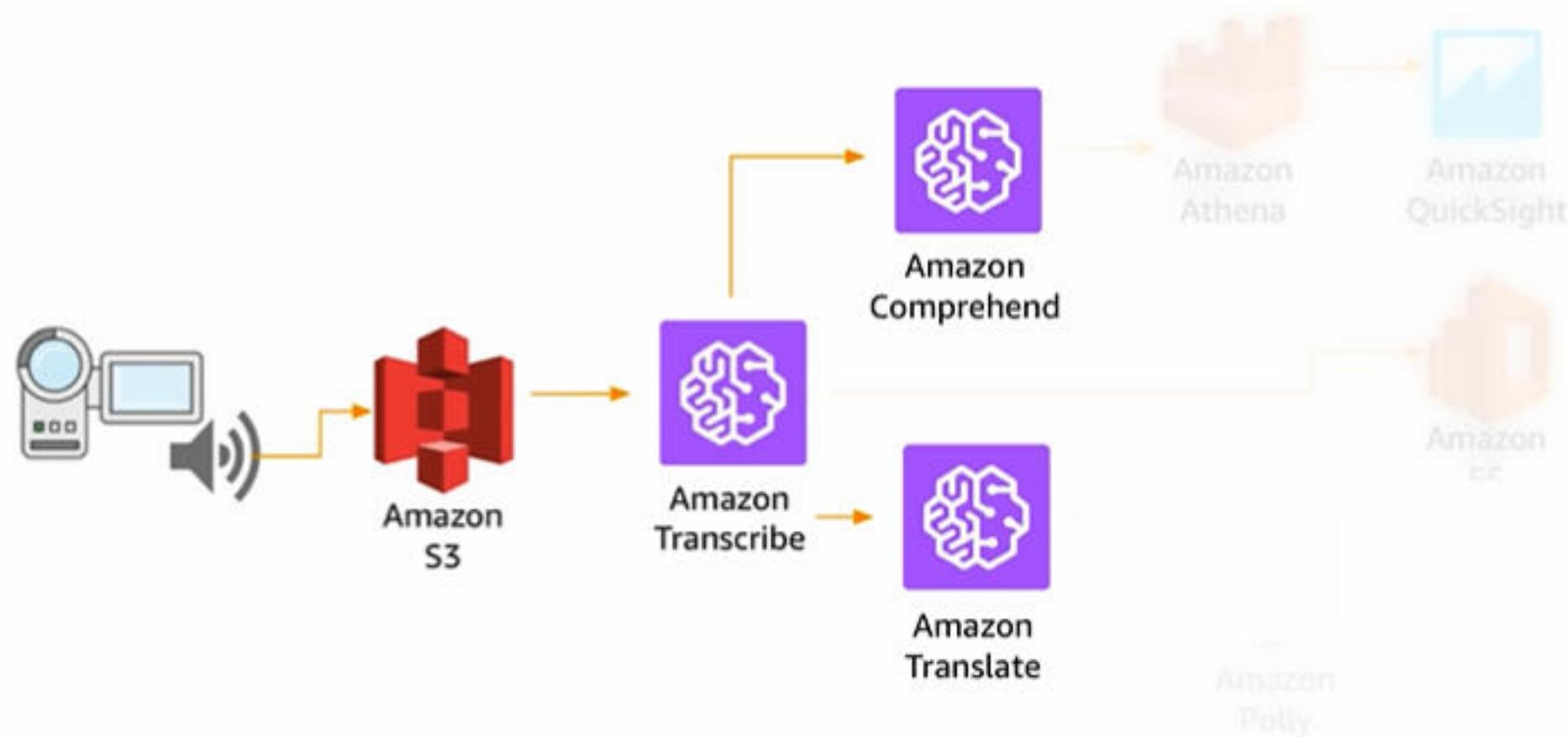
# Use Case Call Center Analysis



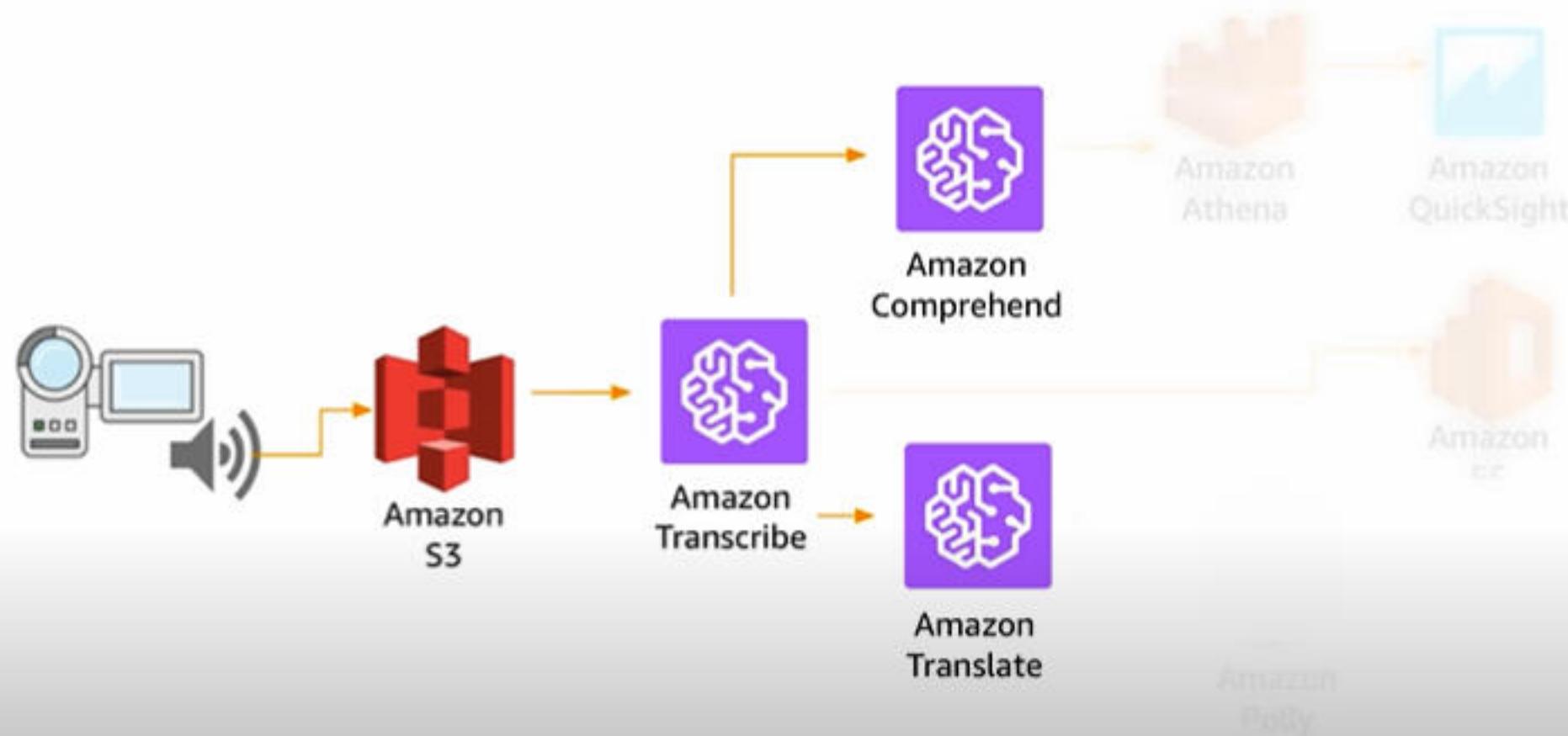
# Use Case: Meeting Transcription



# Use Case: Closed Captioned Video



# Use Case: Closed Captioned Video

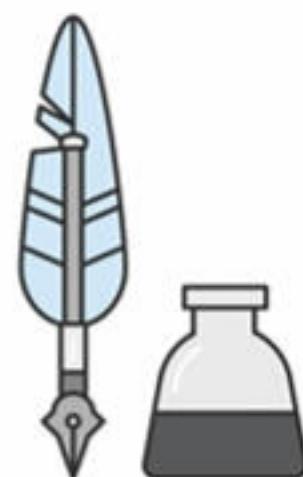


# Getting Started



To try Amazon Transcribe, you'll need:

1. WAV, FLAC, MP3 or MP4 audio file(s) on Amazon S3
2. AWS account with an IAM user that has full access to the Transcribe API calls
3. Familiarity with the Command Line Interface (CLI) and a Text Editor



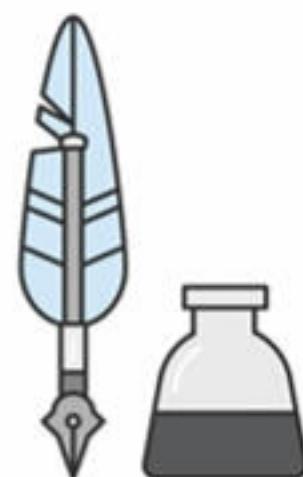
Visit <https://aws.amazon.com/documentation/transcribe/> for the Developer Guide.

# Getting Started



To try Amazon Transcribe, you'll need:

1. WAV, FLAC, MP3 or MP4 audio file(s) on Amazon S3
2. AWS account with an IAM user that has full access to the Transcribe API calls
3. Familiarity with the Command Line Interface (CLI) and a Text Editor



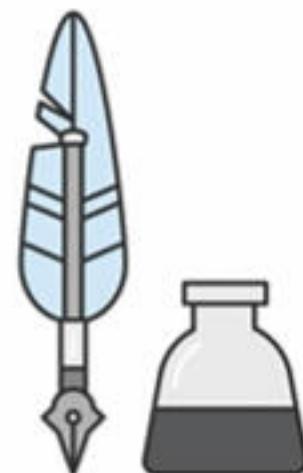
Visit <https://aws.amazon.com/documentation/transcribe/> for the Developer Guide.

# Getting Started



To try Amazon Transcribe, you'll need:

1. WAV, FLAC, MP3 or MP4 audio file(s) on Amazon S3
2. AWS account with an IAM user that has full access to the Transcribe API calls
3. Familiarity with the Command Line Interface (CLI) and a Text Editor



Visit <https://aws.amazon.com/documentation/transcribe/> for the Developer Guide.

# Demonstration



- 큐 Create a transcription job
- 큐 Verify the job's completion
- 큐 Review the results

The screenshot shows the AWS console interface for the Amazon Transcribe service. The top navigation bar includes links for 'Getting Started', 'Corp.', 'TAC', 'Services' (selected), 'Resource Groups', 'Support', and a search bar. The main content area has a dark header with 'MACHINE LEARNING' and the 'Amazon Transcribe' logo. Below the header, the title 'Amazon Transcribe' is displayed in large bold letters, followed by the subtitle 'Automatic Speech Recognition powered by deep learning'. A subtext states: 'Amazon Transcribe provides high-quality and affordable speech-to-text transcription for a wide range of use cases.' To the right, there's a 'Launch the demo' button and a callout box with the text: 'We know you're curious. Jump in and test our Automatic Speech Recognition API.' Below this is a 'Try Amazon Transcribe' button. On the left, under the 'Features' section, there are four items: 'Amazon S3 integration', 'Support for telephony audio', 'Timestamp generation', and 'Easy-to-read transcriptions'. In the bottom left, there's a 'Sample use cases' section with four categories: 'Voice analytics', 'Media & Entertainment', 'Advertising', and 'Search & Compliance'. To the right, there's a 'Pricing (US)' section with a 'Preview' tab showing 'Free' and an 'After preview' section detailing pricing for 1 second of audio at \$0.0004/second. The footer contains links for 'Feedback', 'English (US)', 'Privacy Policy', and 'Terms of Use'.

# Amazon Transcribe

## Automatic Speech Recognition powered by deep learning

Amazon Transcribe provides high-quality and affordable speech-to-text transcription for a wide range of use cases.

### Features

- Amazon S3 integration**  
Amazon Transcribe is integrated with Amazon Simple Storage Service (S3), enabling you to easily transcribe audio files stored in S3 buckets.
- Support for telephony audio**  
Amazon Transcribe supports telephony-quality audio (8 kHz) with high accuracy.
- Timestamp generation**  
Amazon Transcribe returns a timestamp for each word, enabling you to search for a particular word in the transcription and locate it in the original stored recording.
- Easy-to-read transcriptions**  
Amazon Transcribe adds punctuation and formatting automatically so the output is more reader-friendly and can be used without further editing.

### Sample use cases

- Voice analytics**  
Discover customer insights and improve contact center performance.
- Media & Entertainment**  
Automate closed captioning and subtitling workflows for greater accessibility.
- Advertising**
- Search & Compliance**

### Pricing (US)

**Preview**      Free

**After preview**

The first 60 minutes of audio in each monthly cycle will be free for the first year upon signup. Additional usage will be billed at \$0.0004/second on a pay-as-you-go basis. Every request is measured in seconds with a minimum of 15 seconds per request. For additional pricing details, see the pricing page on our website.

1 second of audio	\$0.0004
-------------------	----------

The screenshot shows the 'Create transcription job' page in the Amazon Transcribe service. The page has a dark header bar with the AWS logo and navigation links like 'Services', 'Resource Groups', 'Support', and a search bar. Below the header is a breadcrumb trail: 'Amazon Transcribe > Create transcription job'. The main content area is titled 'Create transcription job' and contains several input fields:

- Input**:
  - Transcription job name**: A text input field labeled 'Enter job name' with placeholder text 'Enter a name for your audio content job.' Below it is a note: 'Job names should be 1-200 characters long, consisting of uppercase and lowercase letters, numbers, and hyphens.'
  - S3 Input URL**: A text input field labeled 'Enter audio content location' with placeholder text 'Paste an S3 link URL to your audio content.' Below it is a note: 'Supported formats: mp3, mp4, wav, flac.'
  - Language**: A dropdown menu set to 'English'.
  - Format**:
    - Format**: A dropdown menu set to 'mp3'.
    - Media sampling rate (Hz) - optional**: An input field labeled 'Enter sampling rate' with placeholder text 'Specify the number of audio samples per second taken per second.' Below it is a note: 'Must be an integer between 8000 and 48000.'

At the bottom right are 'Cancel' and 'Create' buttons. The footer includes links for 'Feedback', 'English (US)', 'Privacy Policy', and 'Terms of Use'.

Concur Hotel | Travel Data... Amazon Transcribe https://console.aws.amazon.com/transcribe/home?region=us-east-1#CreateJob

Most Visited Getting Started Corp TSC

AWS Services Resource Groups Support

Amazon Transcribe > Create transcription job

## Create transcription job

**Input** Info

Transcription job name  
Enter a name for your audio content job.

Job names should be 1-200 characters long, consisting of uppercase and lowercase letters, numbers, and hyphens.

S3 Input URL  
Paste an S3 link URL to your audio content.

Supported formats: mp3, mp4, wav, flac

Language  
Specify the language.

Format Info  
Specify the audio file format that you are uploading.

Supported formats: mp3, mp4, wav, flac

Media sampling rate (Hz) - optional Info  
Specify the number of audio samples per second taken per second.

Must be an integer between 8000 and 48000

Cancel **Create**

Feedback English (US)

© 2006 - 2018, Amazon Web Services, Inc. or its affiliates. All rights reserved. Privacy Policy Terms of Use

**Media Sampling Rate** X

The specified sampling rate (in Hertz) of your audio must match the sampling rate of the associated media content. Sampling rates between 8000 Hz and 48000 Hz are supported by Amazon Transcribe.

Concur Hotel | Travel Data... Amazon Transcribe https://console.aws.amazon.com/transcribe/home?region=us-east-1#jobs

Most Visited Getting Started Corp TSC

AWS Services Resource Groups

Amazon Transcribe

Amazon Transcribe > Transcription jobs

Transcription jobs

Status: All Create job

< 1 2 >

<input type="checkbox"/>	Job name	Language	Creation Time	Status
<input type="checkbox"/>	GettysburgAddressTest	English	Tue Apr 03 2018 15:59:43 GMT-0400 (EDT)	In progress
<input type="checkbox"/>	GettysburgAddress_Test	English	Tue Apr 03 2018 08:10:29 GMT-0400 (EDT)	Complete
<input type="checkbox"/>	TESTGettysburg_Address	English	Mon Apr 02 2018 16:00:24 GMT-0400 (EDT)	Complete
<input type="checkbox"/>	TEST_GettysburgAddress	English	Mon Apr 02 2018 15:56:05 GMT-0400 (EDT)	Complete
<input type="checkbox"/>	TedTalk_2	English	Fri Mar 30 2018 12:13:42 GMT-0400 (EDT)	Complete
<input type="checkbox"/>	GettysburgAddr_2	English	Fri Mar 30 2018 12:12:26 GMT-0400 (EDT)	Complete
<input type="checkbox"/>	Constitution_1	English	Thu Mar 29 2018 19:37:13 GMT-0400 (EDT)	Complete
<input type="checkbox"/>	ARC201	English	Mon Mar 19 2018 19:11:34 GMT-0400 (EDT)	Complete
<input type="checkbox"/>	ML1	English	Mon Mar 05 2018 15:34:10 GMT-0500 (EST)	Complete
<input type="checkbox"/>	Flow	English	Sat Mar 03 2018 20:52:09 GMT-0500 (EST)	Complete

Feedback English (US) © 2006 - 2018, Amazon Web Services, Inc. or its affiliates. All rights reserved. Privacy Policy Terms of Use

S Concur: Hotel | Travel Data... | Amazon Transcribe https://console.aws.amazon.com/transcribe/home?region=us-east-1#jobs

Most Visited - Getting Started Corp TSC

AWS Services Resource Groups

Amazon Transcribe

Amazon Transcribe > Transcription jobs

Transcription jobs

Status: All Download Create job

< 1 2 > ⌂

<input type="checkbox"/>	Job name	Language	Creation Time	Status
<input type="checkbox"/>	GettysburgAddressTest	English	Tue Apr 03 2018 15:59:43 GMT-0400 (EDT)	<span>In progress</span>
<input type="checkbox"/>	GettysburgAddress_Test	English	Tue Apr 03 2018 08:10:29 GMT-0400 (EDT)	<span>Complete</span>
<input type="checkbox"/>	TESTGettysburg_Address	English	Mon Apr 02 2018 16:00:24 GMT-0400 (EDT)	<span>Complete</span>
<input type="checkbox"/>	TEST_GettysburgAddress	English	Mon Apr 02 2018 15:56:05 GMT-0400 (EDT)	<span>Complete</span>
<input type="checkbox"/>	TedTalk_2	English	Fri Mar 30 2018 12:13:42 GMT-0400 (EDT)	<span>Complete</span>
<input type="checkbox"/>	GettysburgAddr_2	English	Fri Mar 30 2018 12:12:26 GMT-0400 (EDT)	<span>Complete</span>
<input type="checkbox"/>	Constitution_1	English	Thu Mar 29 2018 19:37:13 GMT-0400 (EDT)	<span>Complete</span>
<input type="checkbox"/>	ARC201	English	Mon Mar 19 2018 19:11:34 GMT-0400 (EDT)	<span>Complete</span>
<input type="checkbox"/>	ML1	English	Mon Mar 05 2018 15:34:10 GMT-0500 (EST)	<span>Complete</span>
<input type="checkbox"/>	Flow	English	Sat Mar 03 2018 20:52:09 GMT-0500 (EST)	<span>Complete</span>

Feedback English (US) © 2006 - 2018, Amazon Web Services, Inc. or its affiliates. All rights reserved. Privacy Policy Terms of Use

Warning Amazon Transcribe https://s3.amazonaws.com/GettysburgAddress.mp3

https://console.aws.amazon.com/transcribe/home?region=us-east-1#job-details/GettysburgAddressTest

Most Visited Getting Started Corp T&C

aWS Services Resource Groups Support

Amazon Transcribe > Transcription jobs > GettysburgAddressTest

## GettysburgAddressTest

**Details**

Job name  
GettysburgAddressTest

Status  
 Complete

Expiration date  
The transcription will be available for another 89 days.

Input URL  
[https://s3.amazonaws.com/currdev-transcribe/Gettysburg\\_Address.mp3](https://s3.amazonaws.com/currdev-transcribe/Gettysburg_Address.mp3)

Output URL  
<https://s3.amazonaws.com/> /GettysburgAddressTest/asrOutput.json

**Download transcription**

**Code Samples**

**Audio conversion**

JSON Request

```
{ "TranscriptionJobName": "GettysburgAddressTest", "LanguageCode": "en-US", "MediaSampleRateHertz": 44100, "MediaFormat": "mp3", "Media": { "MediaFileUrl": "https://s3.amazonaws.com/currdev-transcribe/Gettysburg_Address.mp3" } }
```

**Select**

**View API documentation**

Feedback English (US) Privacy Policy Terms of Use

|| 10 08:43 / 09:41

https://www...essing%25.asp

Amazon Transcribe https://s3.amazonaws.com/GettysburgAddress.mp3

https://console.aws.amazon.com/transcribe/home?region=us-east-1#job-details/GettysburgAddressTest

Most Visited Getting Started Corp T&C

aWS Services Resource Groups

Amazon Transcribe > Transcription jobs > GettysburgAddressTest

## GettysburgAddressTest

**Details**

Job name  
GettysburgAddressTest

Status  
 Complete

Expiration date  
The transcription will be available for another 89 days.

Input URL  
[https://s3.amazonaws.com/currdev-transcribe/Gettysburg\\_Address.mp3](https://s3.amazonaws.com/currdev-transcribe/Gettysburg_Address.mp3)

Output URL  
<https://s3.amazonaws.com/GettysburgAddressTest/asrOutput.json>

**Code Samples**

View API documentation

**Audio conversion**

JSON Request

```
{ "TranscriptionJobName": "GettysburgAddressTest", "LanguageCode": "en-US", "MediaSampleRateHertz": 44100, "MediaFormat": "mp3", "Media": { "MediaFileUrl": "https://s3.amazonaws.com/currdev-transcribe/Gettysburg_Address.mp3" }}
```

Select

https://s3.amazonaws.com/aws-transcribe-us-east-1-prod/820309686234/GettysburgAddressTest/asrOutput.json?X-Amz-Security-Token=FQoDYXdzEEAjdK+YD7aODh6e6Bsa... X-Amz-Credential=ASIAjT752DXN7UxVR2Q/20180403/us-east-1/s3/aws4\_request&X-Amz-Signature=88375887709da89de3fd7a2c8ed31562a52d7cf6b3c8ac7014249144f883176a

https://www...essing%2fasp... Amazon Transcribe https://s3.amazonaws.com/GettysburgAddressTest.mp3

https://console.aws.amazon.com/transcribe/home?region=us-east-1#job-details/GettysburgAddressTest

Most Visited Getting Started Corp T&C

aWS Services Resource Groups Support

Amazon Transcribe > Transcription jobs > GettysburgAddressTest

## GettysburgAddressTest

**Details**

Job name  
GettysburgAddressTest

Status  
 Complete

Expiration date  
The transcription will be available for another 89 days.

Input URL  
[https://s3.amazonaws.com/currdev-transcribe/Gettysburg\\_Address.mp3](https://s3.amazonaws.com/currdev-transcribe/Gettysburg_Address.mp3)

Output URL  
<https://s3.amazonaws.com/GettysburgAddressTest/asrOutput.json>

**Code Samples**

View API documentation

**Audio conversion**

JSON Request

```
{ "TranscriptionJobName": "GettysburgAddressTest", "LanguageCode": "en-US", "MediaSampleRateHertz": 44100, "MediaFormat": "mp3", "Media": { "MediaFileUrl": "https://s3.amazonaws.com/currdev-transcribe/Gettysburg_Address.mp3" }}
```

Select

https://s3.amazonaws.com/aws-transcribe-us-east-1-prod/820309686234/GettysburgAddressTest/asrOutput.json?X-Amz-Security-Token=FQoDYXdzEEAjdK+YD7aODh6e6Bsa... X-Amz-Credential=ASIAjT752DXN7UxVR2Q/20180403/us-east-1/s3/aws4\_request&X-Amz-Signature=88375887709da89de3fd7a2c8ed31562a52d7cf6b3c8ac7014249144f883176a

TextEdit File Edit Format View Window Help

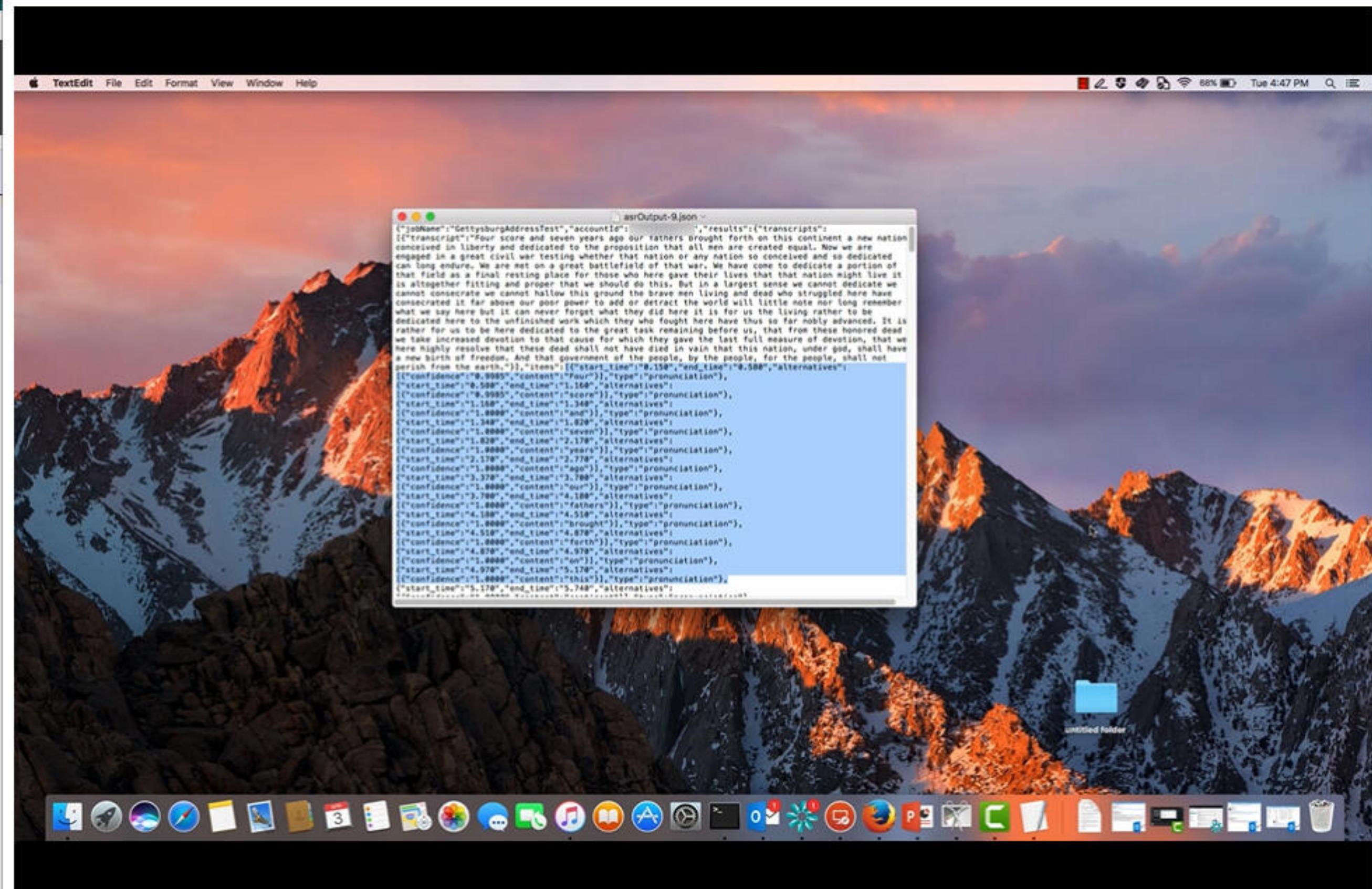
88% Tue 4:47 PM

TextEdit

asrOutput-9.json

```
{"jobName": "GettysburgAddressfest", "accountId": "12345678901234567890123456789012", "results": {"transcripts": [{"transcript": "Four score and seven years ago our fathers brought forth on this continent a new nation conceived in liberty and dedicated to the proposition that all men are created equal. Now we are engaged in a great civil war testing whether that nation or any nation so conceived and so dedicated can long endure. We are met on a great battlefield of that war. We have come to dedicate a portion of that field as a final resting place for those who here gave their lives that that nation might live it is altogether fitting and proper that we should do this. But in a largest sense we cannot dedicate we cannot consecrate we cannot hallow this ground the brave men living and dead who struggled here have consecrated it far above our poor power to add or detract the world will little note nor long remember what we say here but it can never forget what they did here it is for us the living rather to be dedicated here to the unfinished work which they who fought here have thus so far nobly advanced. It is rather for us to be here dedicated to the great task remaining before us, that from these honored dead we take increased devotion to that cause for which they gave the last full measure of devotion, that we here highly resolve that these dead shall not have died in vain that this nation, under god, shall have a new birth of freedom. And that government of the people, by the people, for the people, shall not perish from the earth."}], "items": [{"start_time": "0.150", "end_time": "0.580", "alternatives": [{"confidence": "0.9995", "content": "Four"}]}, {"start_time": "0.580", "end_time": "1.160", "alternatives": [{"confidence": "0.9995", "content": "score"}]}, {"start_time": "1.160", "end_time": "1.340", "alternatives": [{"confidence": "1.0000", "content": "and"}]}, {"start_time": "1.340", "end_time": "1.820", "alternatives": [{"confidence": "1.0000", "content": "ago"}]}, {"start_time": "1.820", "end_time": "2.170", "alternatives": [{"confidence": "1.0000", "content": "years"}]}, {"start_time": "2.170", "end_time": "2.770", "alternatives": [{"confidence": "1.0000", "content": "brought"}]}, {"start_time": "2.770", "end_time": "3.370", "alternatives": [{"confidence": "1.0000", "content": "forth"}]}, {"start_time": "3.370", "end_time": "3.760", "alternatives": [{"confidence": "1.0000", "content": "our"}]}, {"start_time": "3.760", "end_time": "4.160", "alternatives": [{"confidence": "1.0000", "content": "fathers"}]}, {"start_time": "4.160", "end_time": "4.510", "alternatives": [{"confidence": "1.0000", "content": "brought"}]}, {"start_time": "4.510", "end_time": "4.870", "alternatives": [{"confidence": "1.0000", "content": "forth"}]}, {"start_time": "4.870", "end_time": "4.970", "alternatives": [{"confidence": "1.0000", "content": "on"}]}, {"start_time": "4.970", "end_time": "5.170", "alternatives": [{"confidence": "1.0000", "content": "this"}]}, {"start_time": "5.170", "end_time": "5.740", "alternatives": [{"confidence": "1.0000", "content": "will"}]}]}
```







## Thanks for watching!

© 2018 Amazon Web Services, Inc. or its affiliates. All rights reserved. This work may not be reproduced or redistributed, in whole or in part, without prior written permission from Amazon Web Services, Inc. Commercial copying, lending, or selling is prohibited. Corrections or feedback on the course, please email us at: [aws-course-feedback@amazon.com](mailto:aws-course-feedback@amazon.com). For all other questions, contact us at: <https://aws.amazon.com/contact-us/aws-training/>. All trademarks are the property of their owners.

## Knowledge Acquisition



## Inference



# Types of Machine Learning



## Supervised



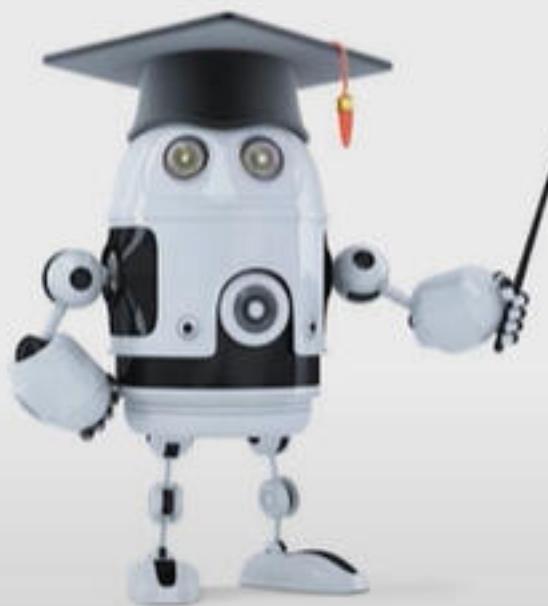
## Unsupervised



# Types of Machine Learning



## Supervised



## Unsupervised

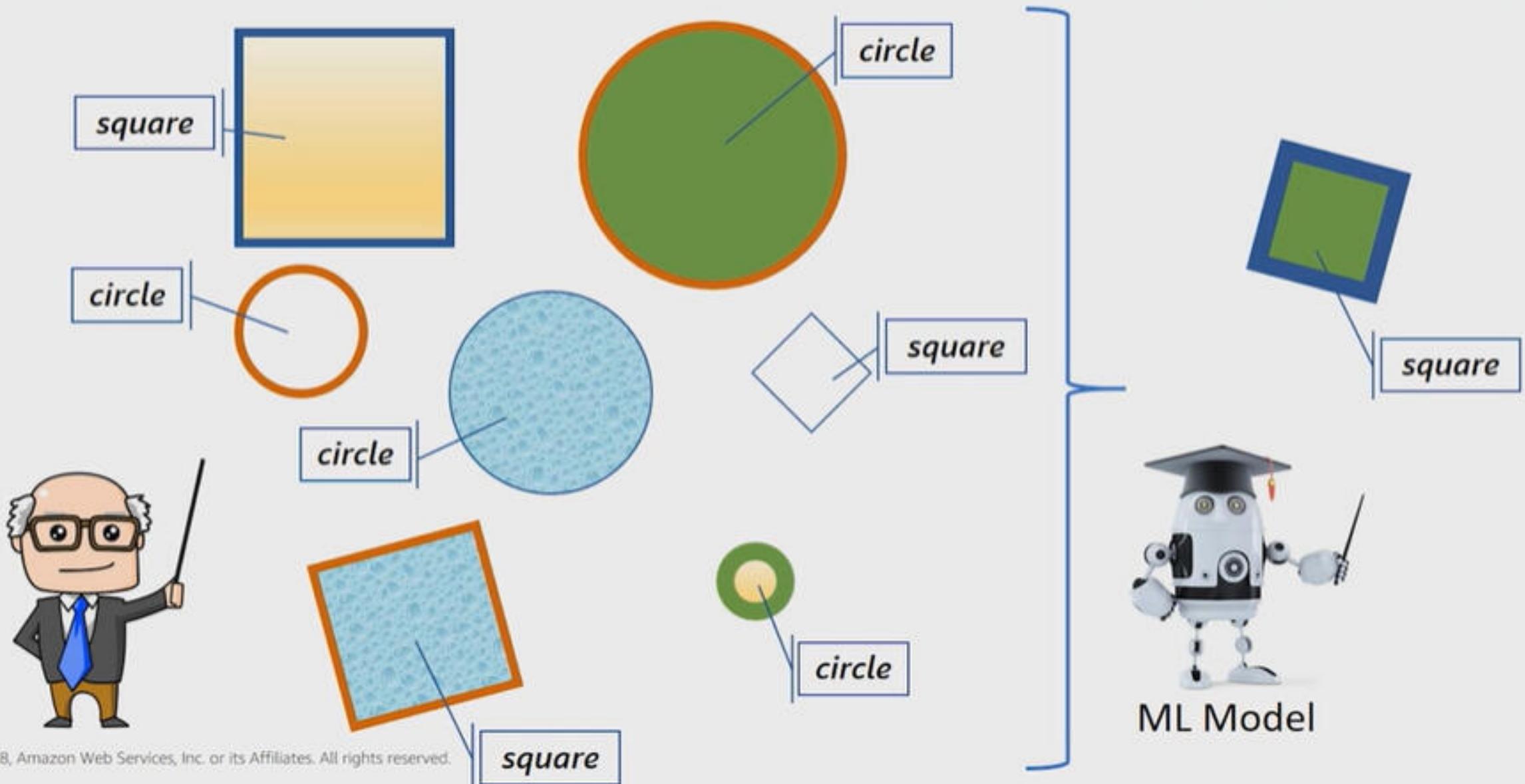
things I can  
eat      play with



## Reinforcement



# Supervised Learning: Example



© 2018, Amazon Web Services, Inc. or its Affiliates. All rights reserved.

# Supervised Learning

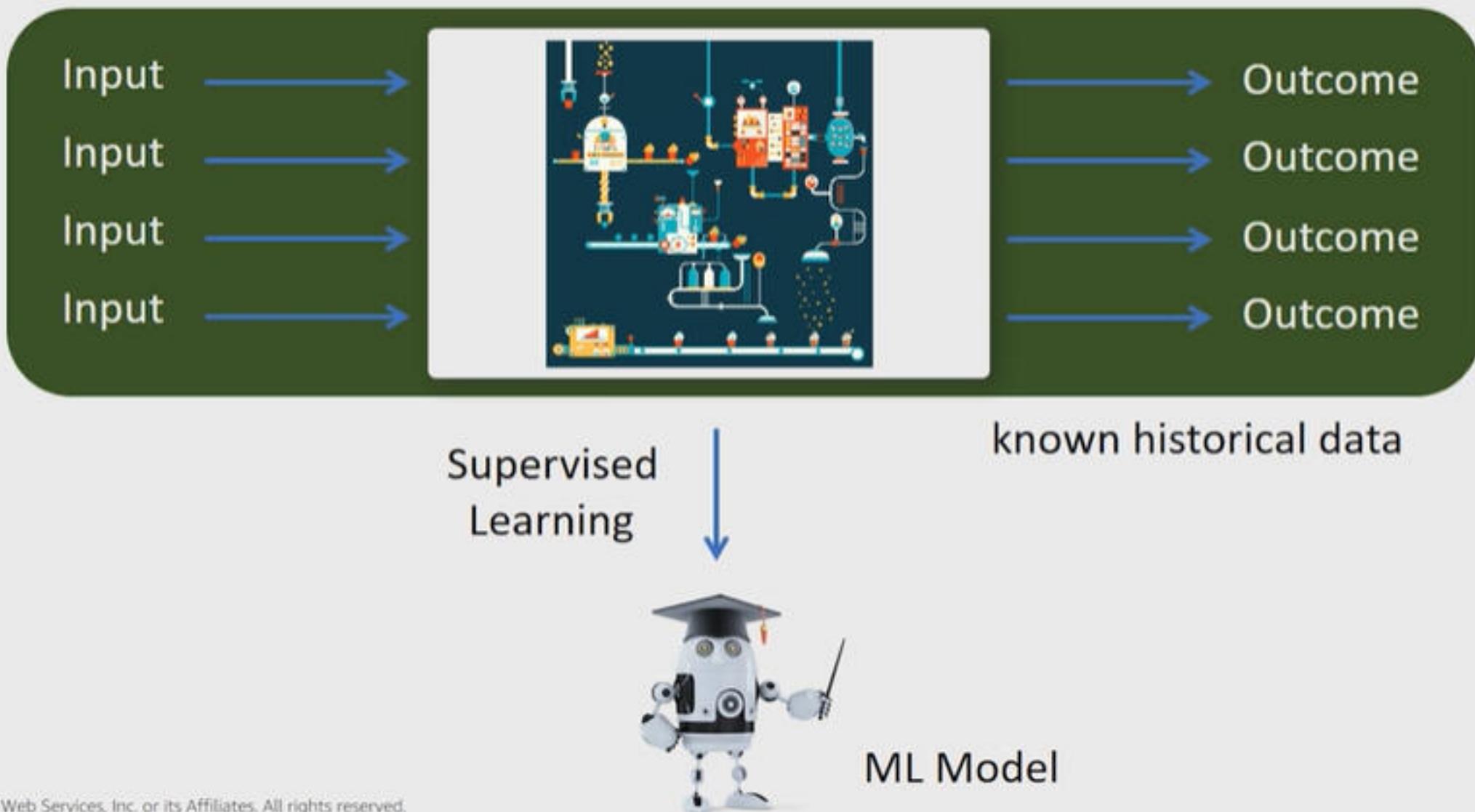
aws training and certification



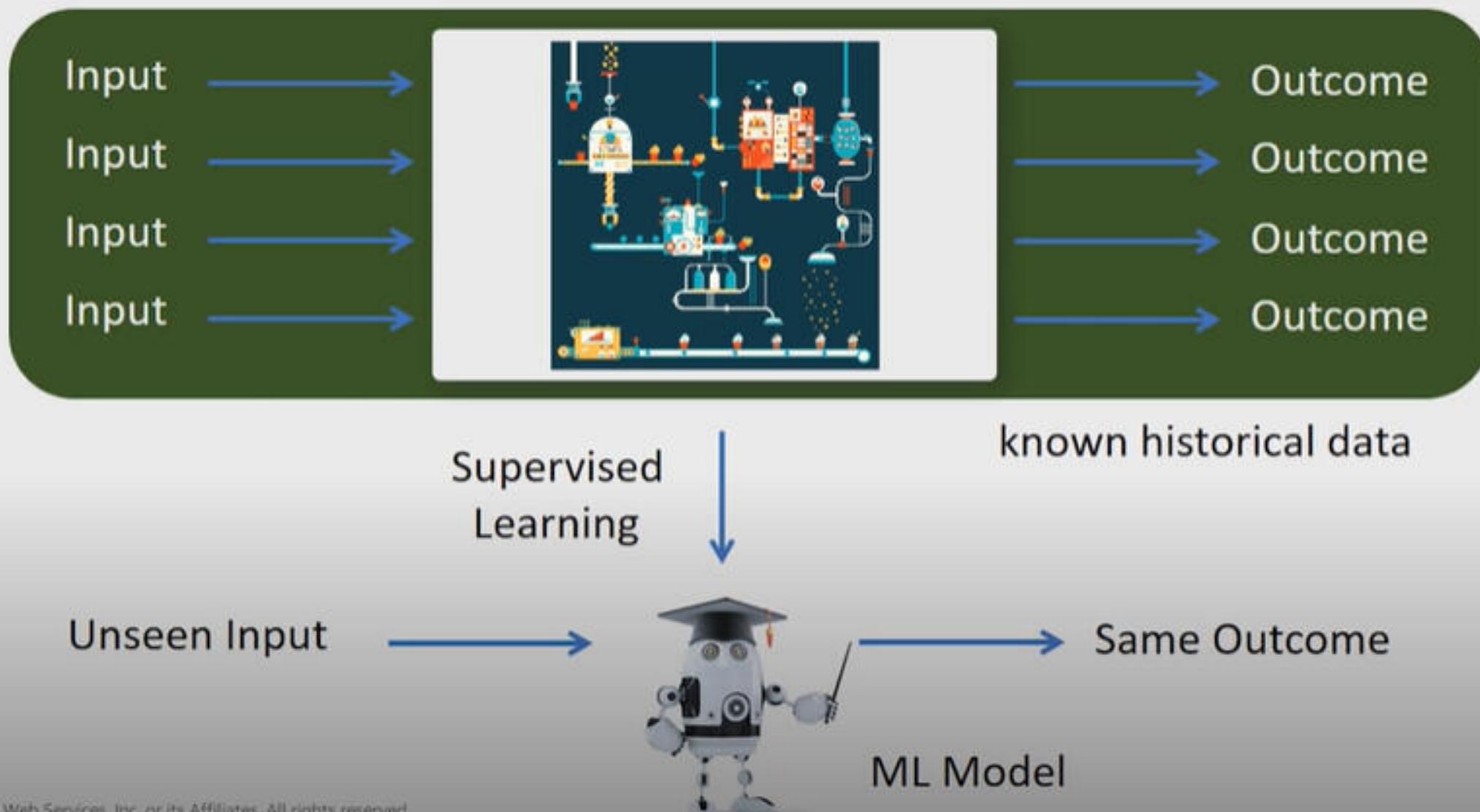
# Supervised Learning: Observations



# Supervised Learning: Training

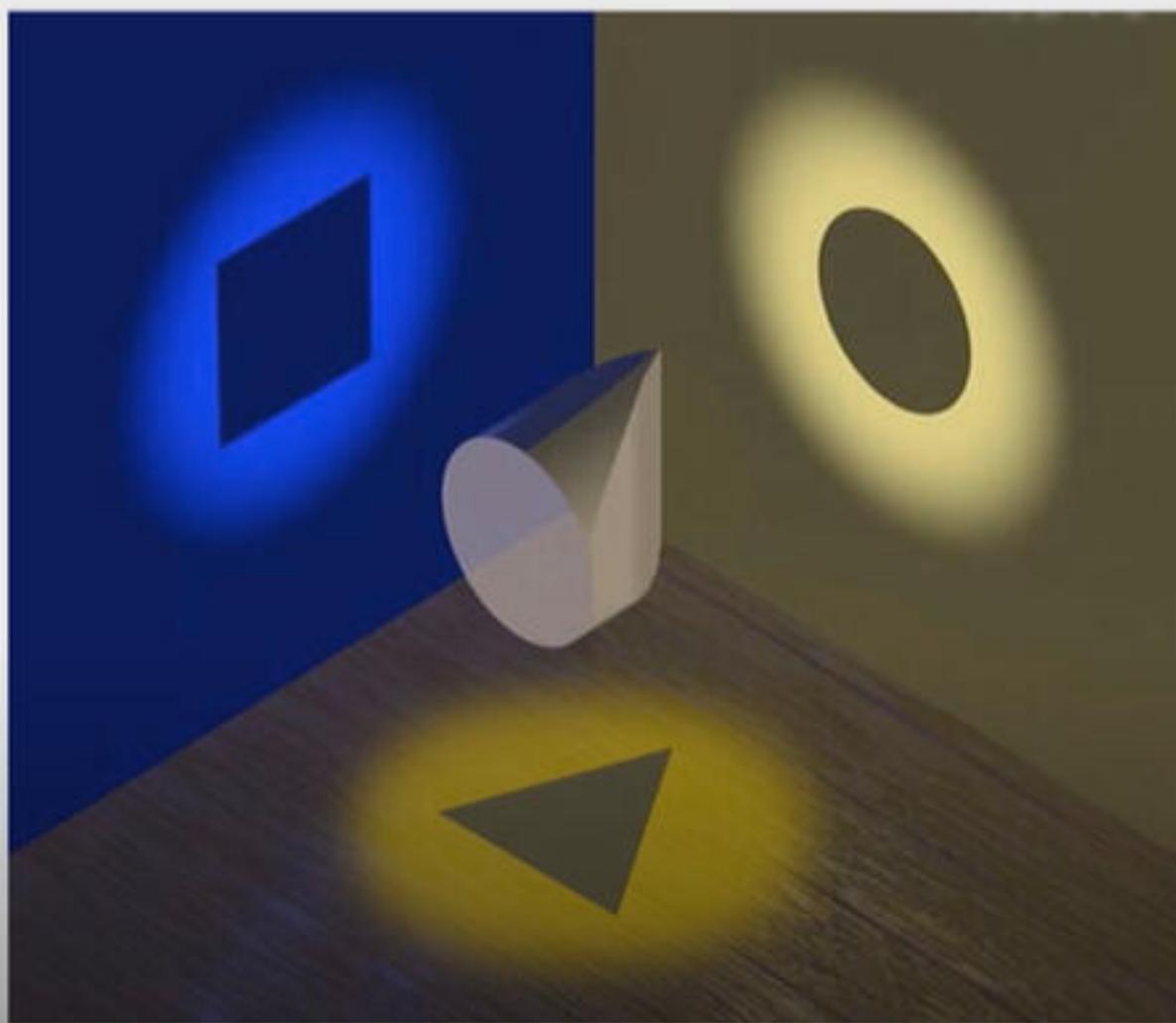


# Supervised Learning: Trained Model

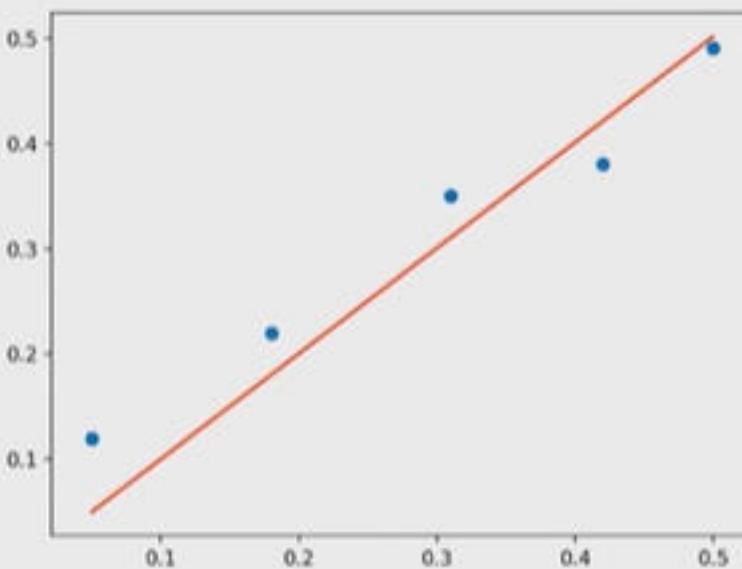


# Mutually Exclusive Classes

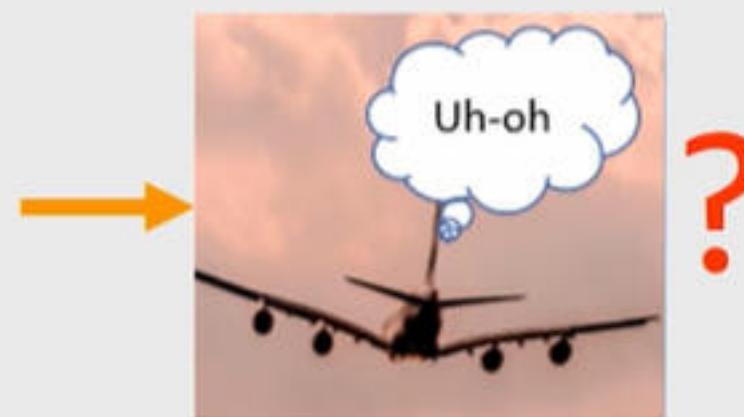
aws training and certification



# Predicting Numbers: Regression



# Supervised Learning Examples



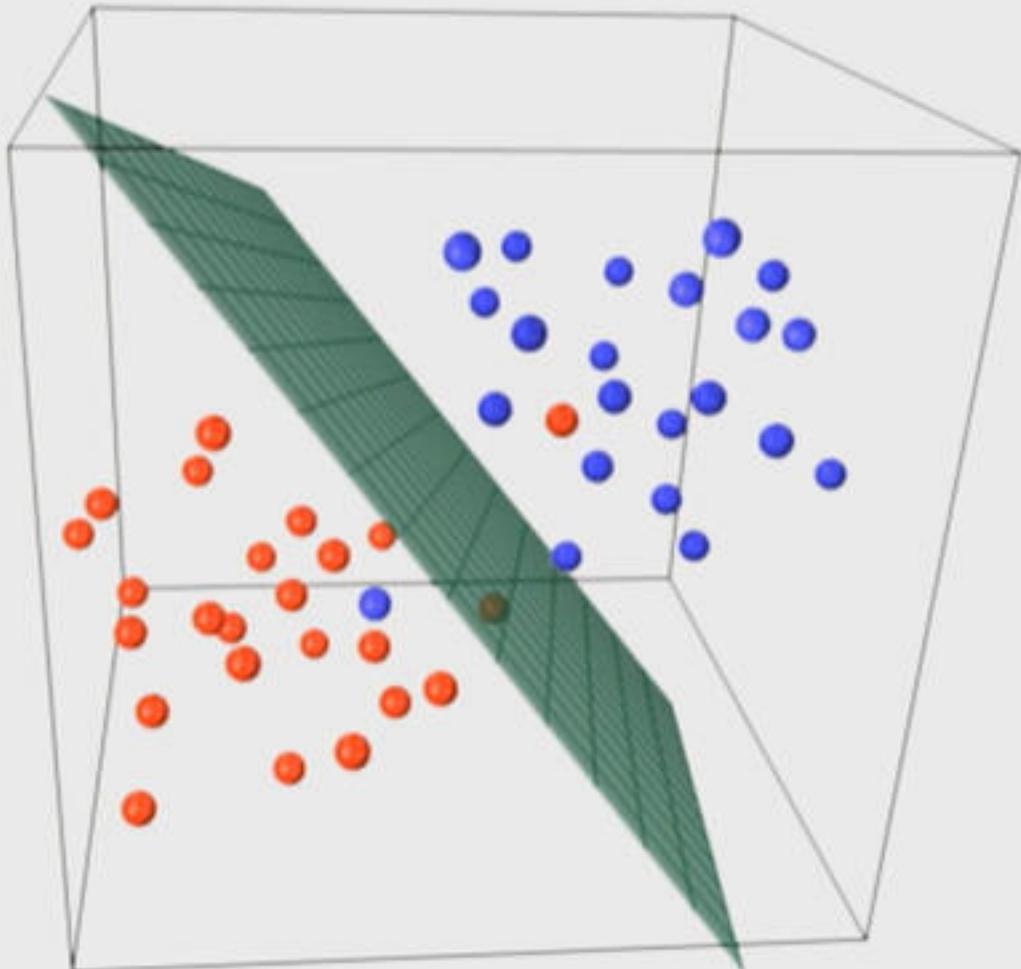
# Supervised Learning Examples



# Supervised Learning Algorithms

# Linear Supervised Algorithms

aws training and certification



# Linear Supervised Algorithms



- In Amazon SageMaker: Linear Learner
  - Linear + Logistic Regression



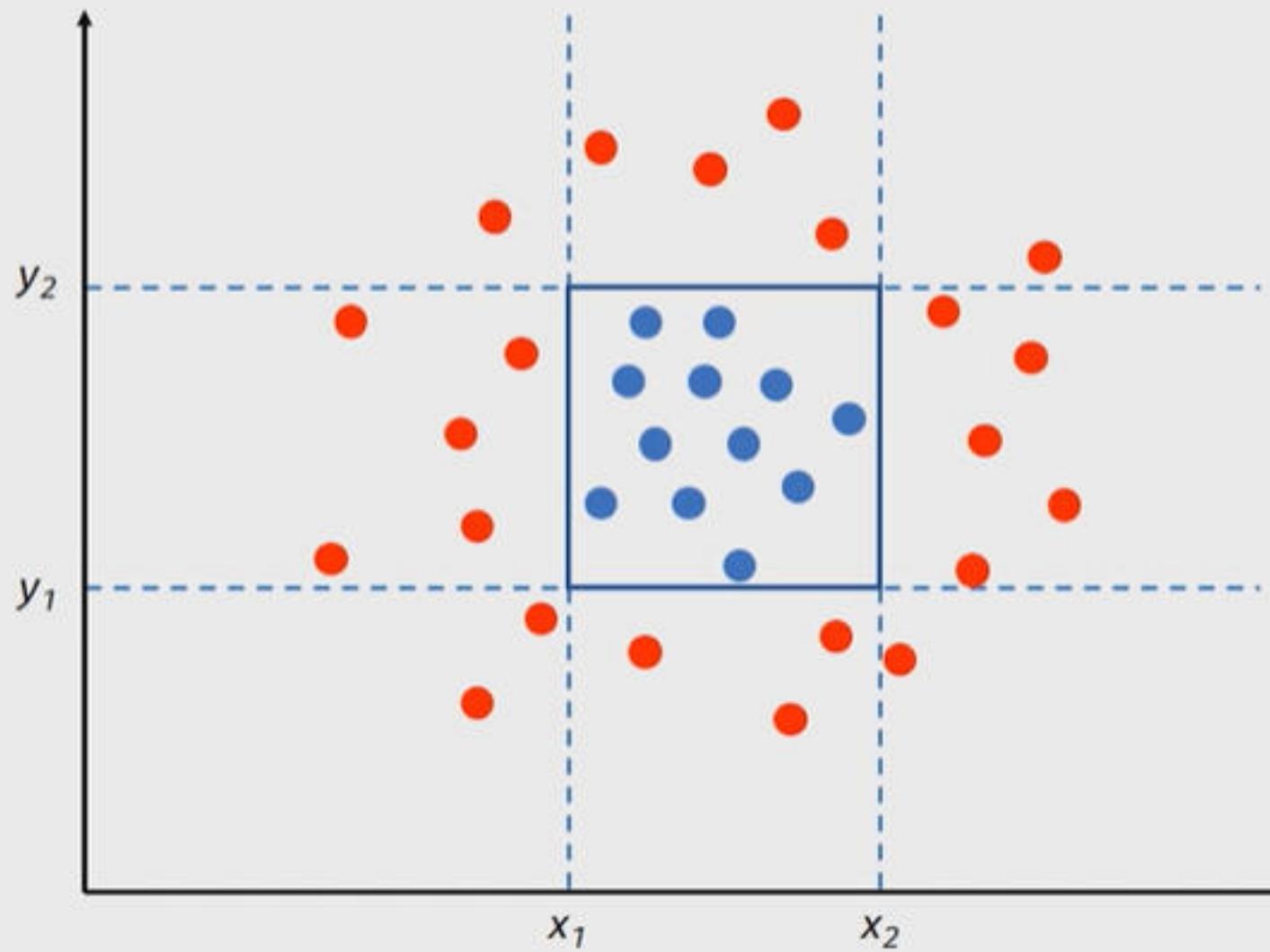
# Linear Supervised Algorithms



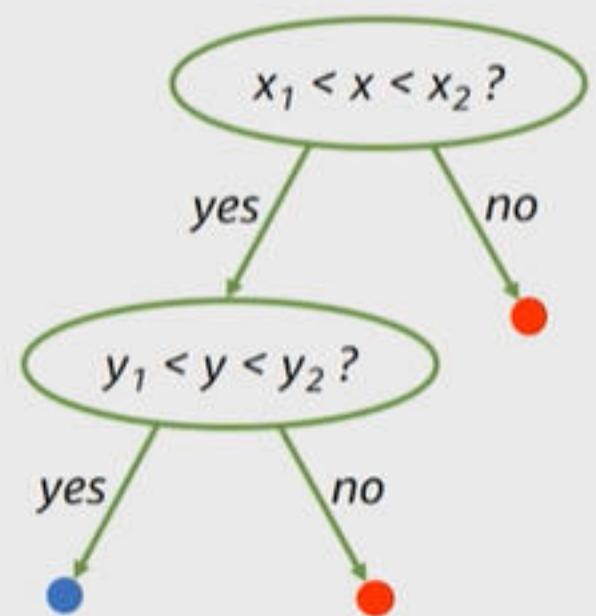
- In Amazon SageMaker: Linear Learner
  - Linear + Logistic Regression
- Other common algorithms:
  - Support Vector Machine (SVM)
  - Perceptron



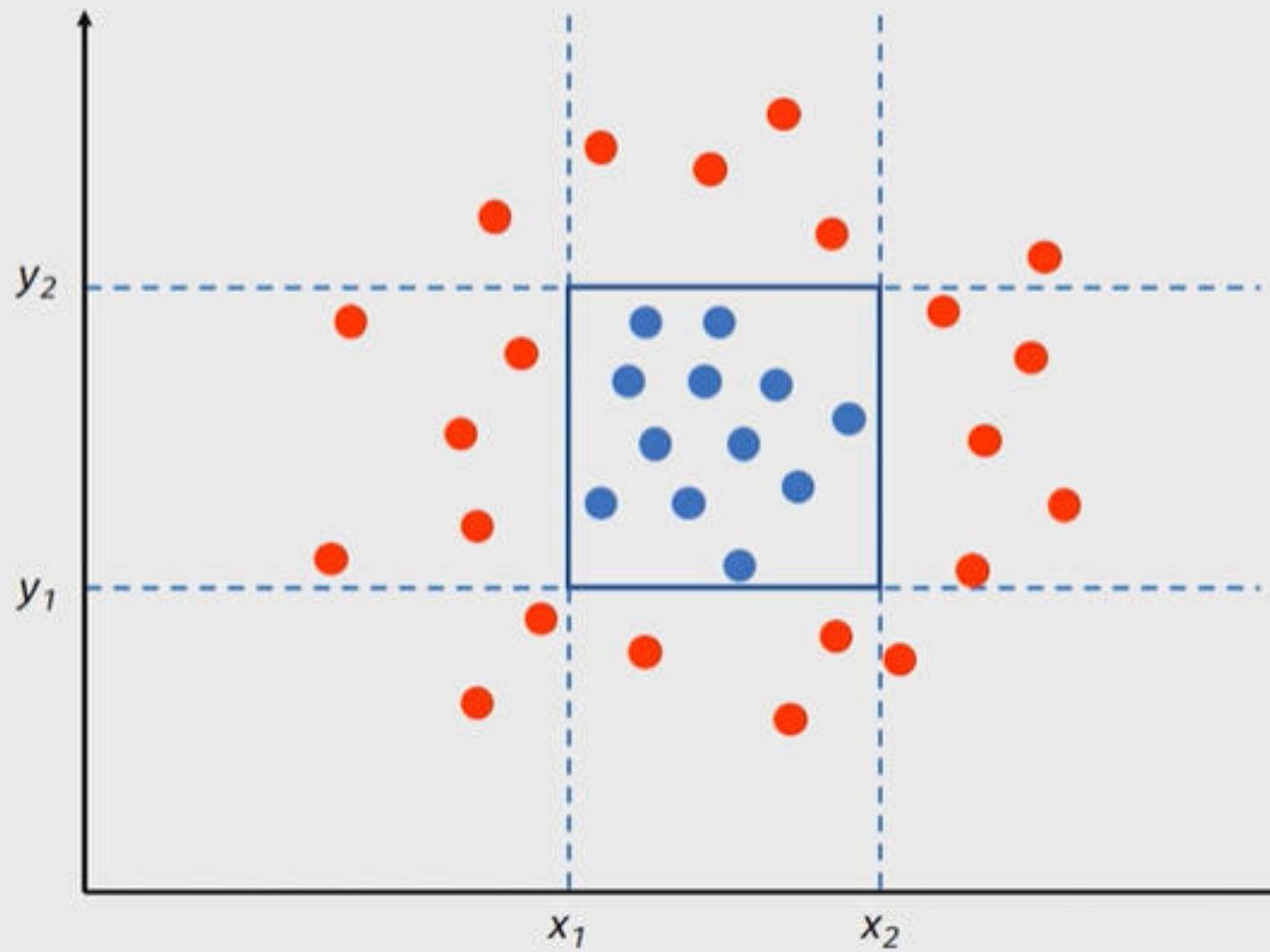
# Non-Linear Supervised Algorithms



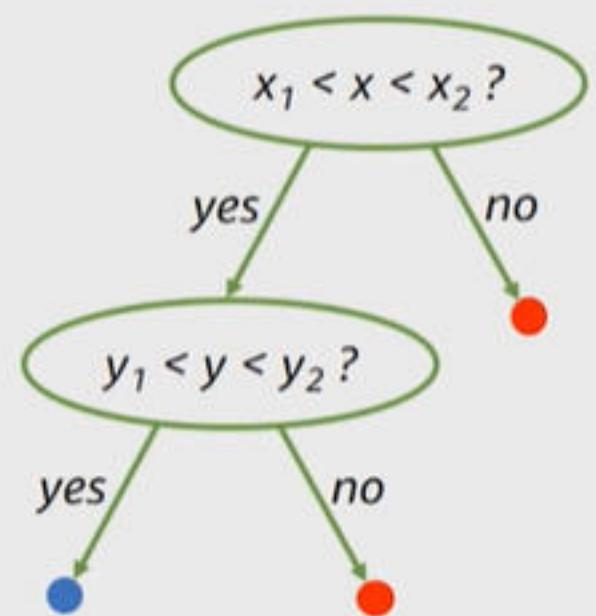
Decision Tree



# Non-Linear Supervised Algorithms



Decision Tree



# Non-Linear Supervised Algorithms



- XGBoost - Gradient Boosted Trees
- Factorization Machines
  - Good for high dimensional sparse datasets
  - e.g. click prediction & item recommendation



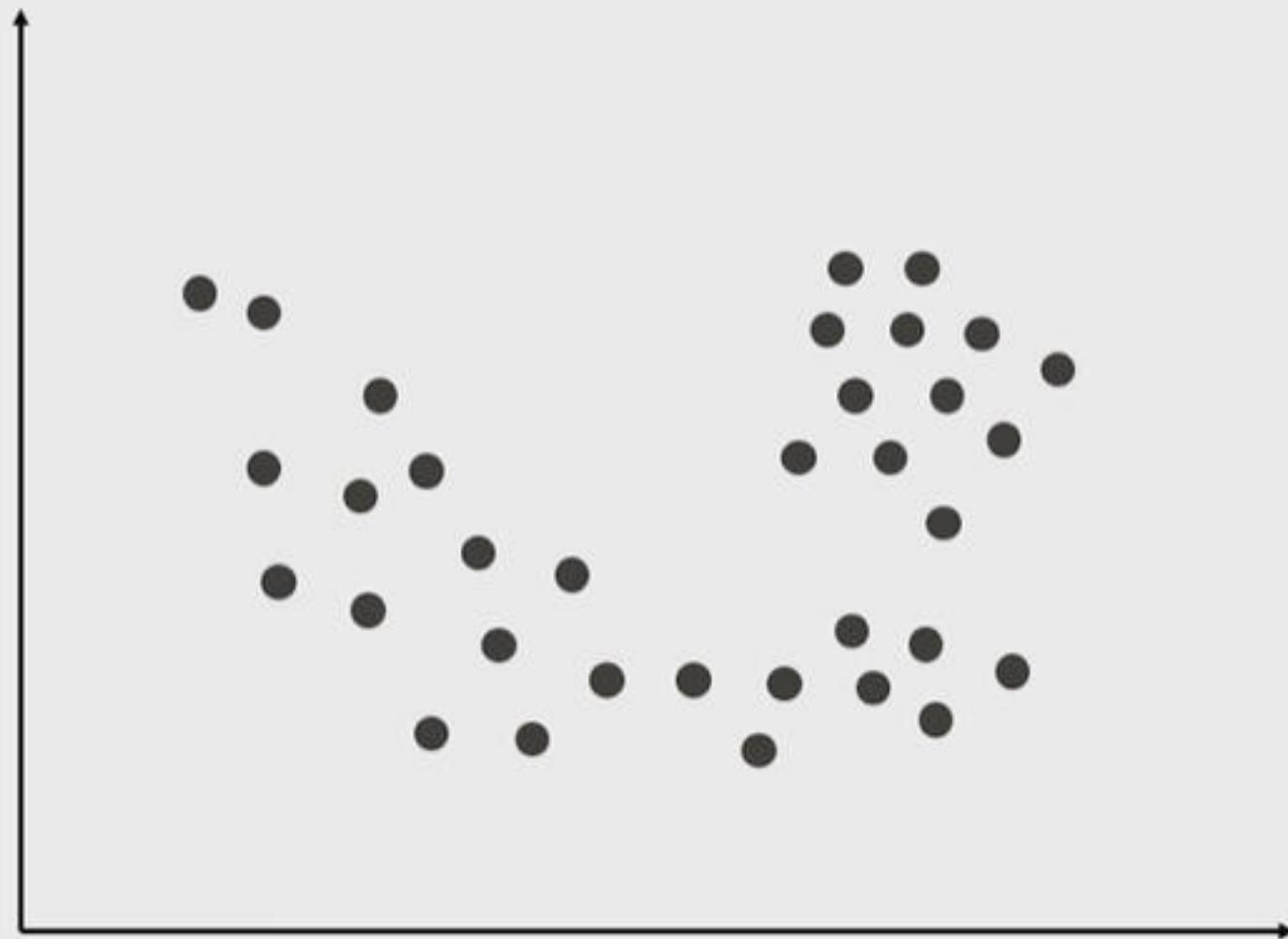
# Non-Linear Supervised Algorithms



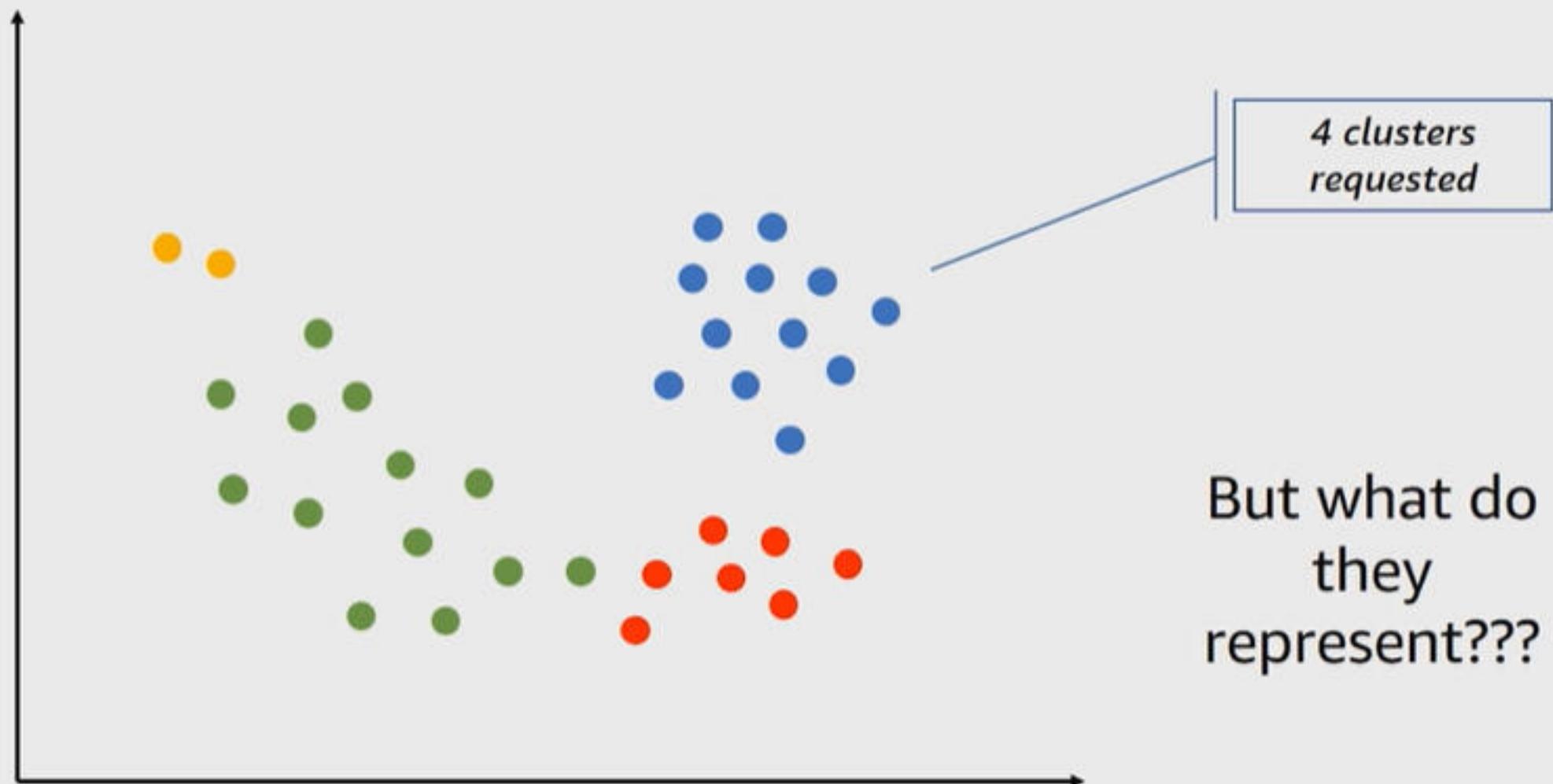
- XGBoost - Gradient Boosted Trees
- Factorization Machines
  - Good for high dimensional sparse datasets
  - e.g. click prediction & item recommendation
- Other popular approaches
  - Polynomial
  - Neural Networks



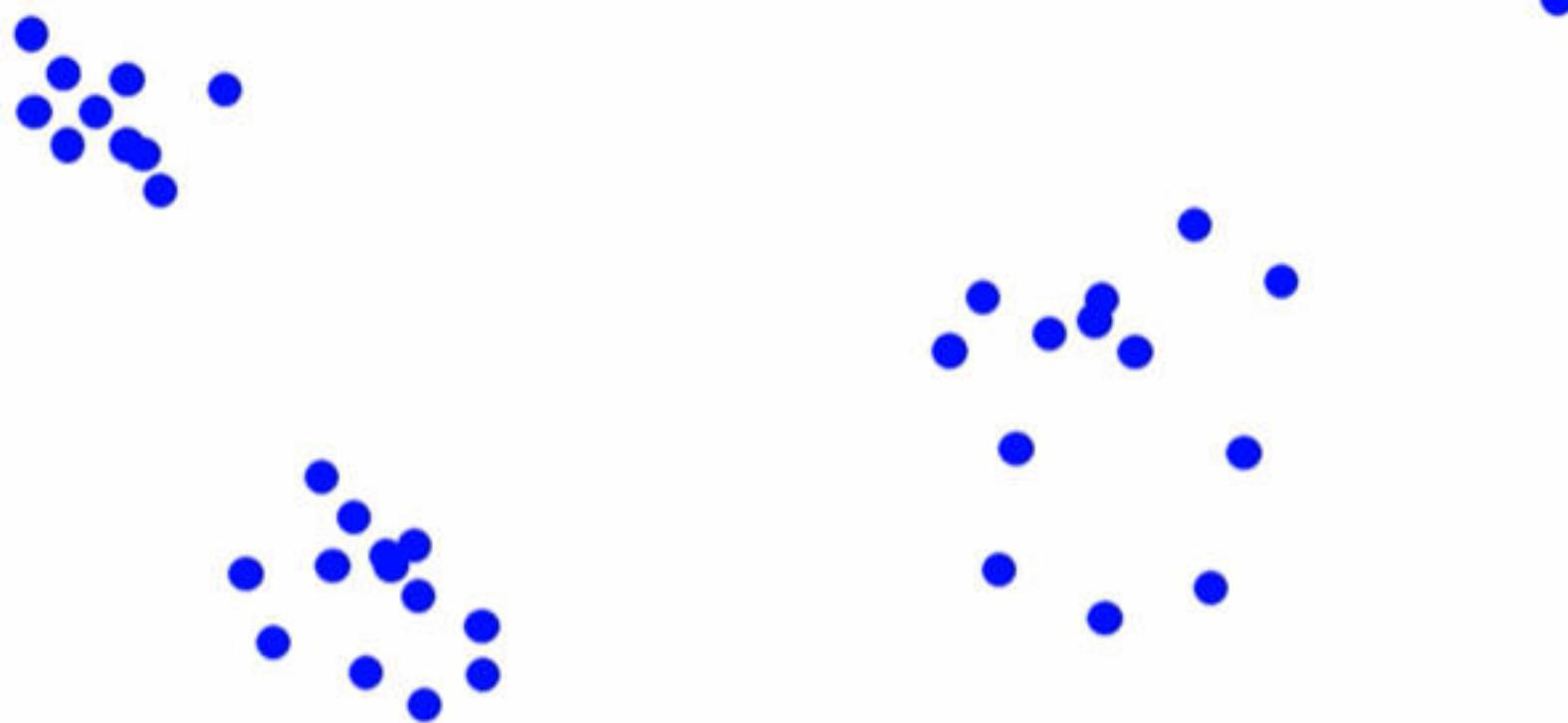
# Unsupervised Learning: Clustering



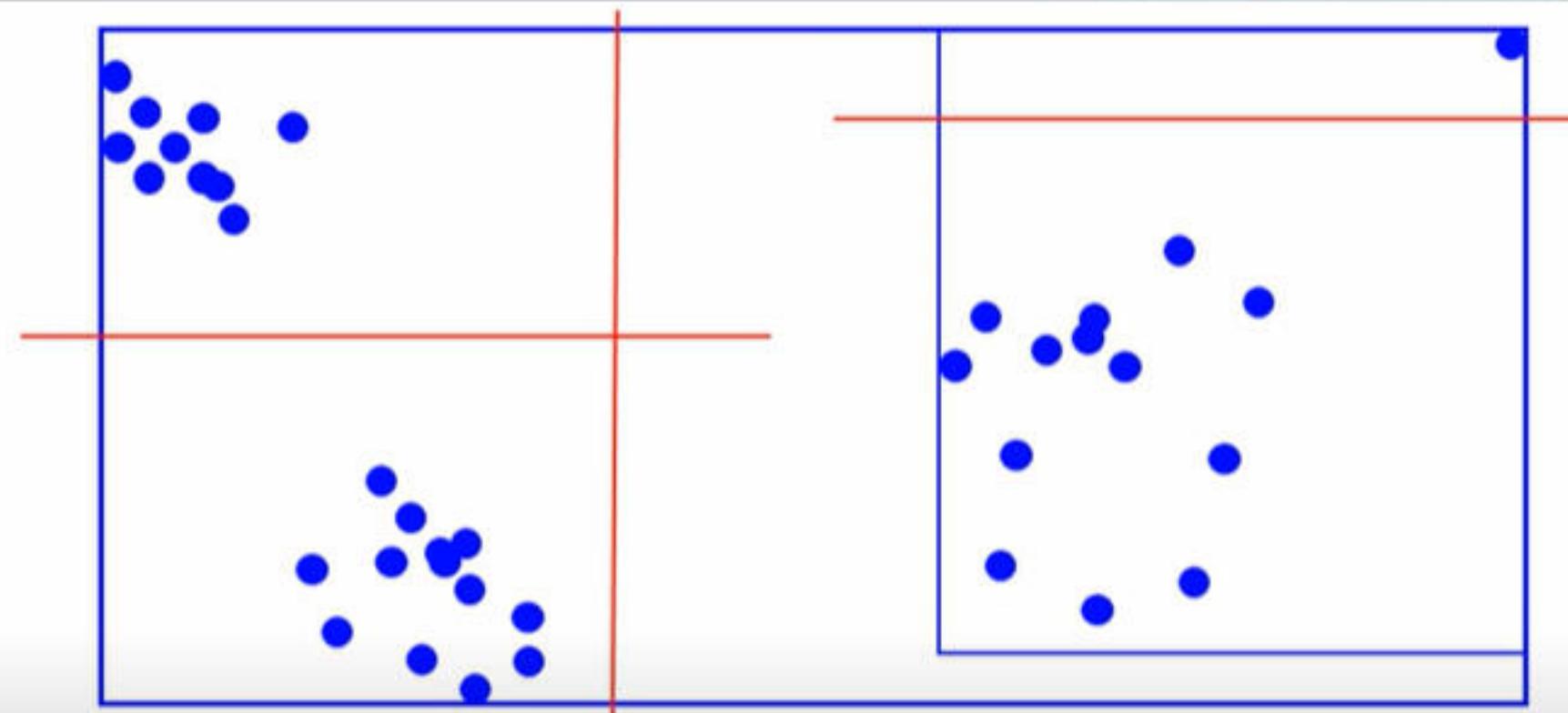
# Unsupervised Learning: Clustering



# Anomalies: Random Cut Forest



# Anomalies: Random Cut Forest



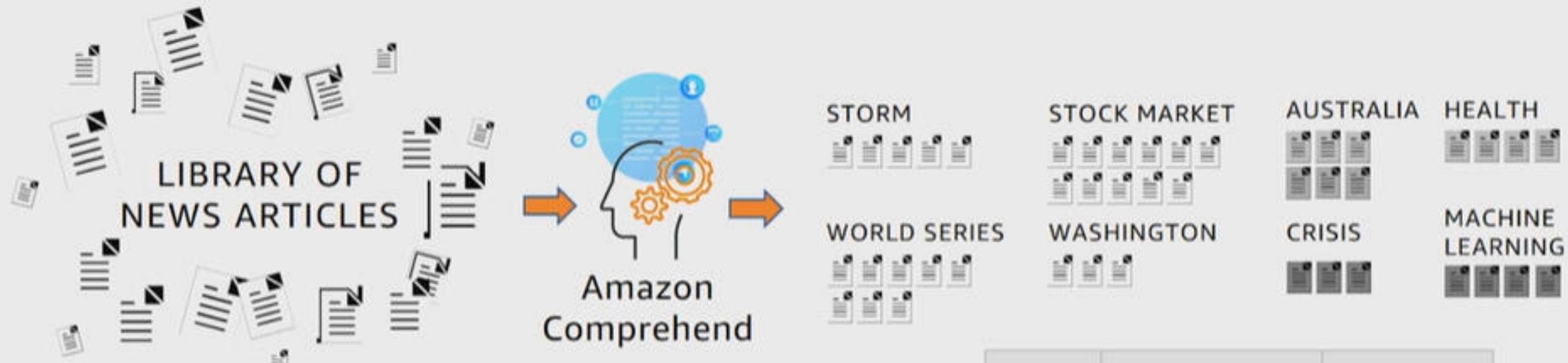
Range-biased Cut

Recurse: The cutting stops when each point is isolated.

# Unsupervised ML: Topic Modeling



# Unsupervised ML: Topic Modeling

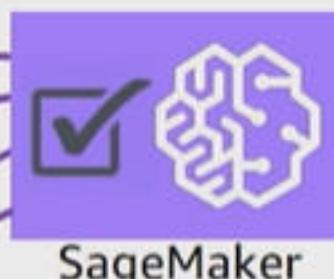


Topic	Term	Weight
000	team	0.118533
000	game	0.106072
000	player	0.031625
000	season	0.023633

# Unsupervised Algorithms



- K-Means Clustering
  - Improved Web-scale k-means clustering
- Principal Component Analysis (PCA)
  - Reduces dimensionality within dataset
  - Precursor to Supervised Learning
- Latent Dirichlet Allocation (LDA)
  - Topic Modelling
- Anomaly Detection



SageMaker



Comprehend

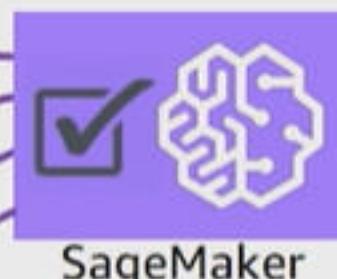


Kinesis Data  
Analytics

# Unsupervised Algorithms



- K-Means Clustering
  - Improved Web-scale k-means clustering
- Principal Component Analysis (PCA)
  - Reduces dimensionality within dataset
  - Precursor to Supervised Learning
- Latent Dirichlet Allocation (LDA)
  - Topic Modelling
- Anomaly Detection



SageMaker



Comprehend

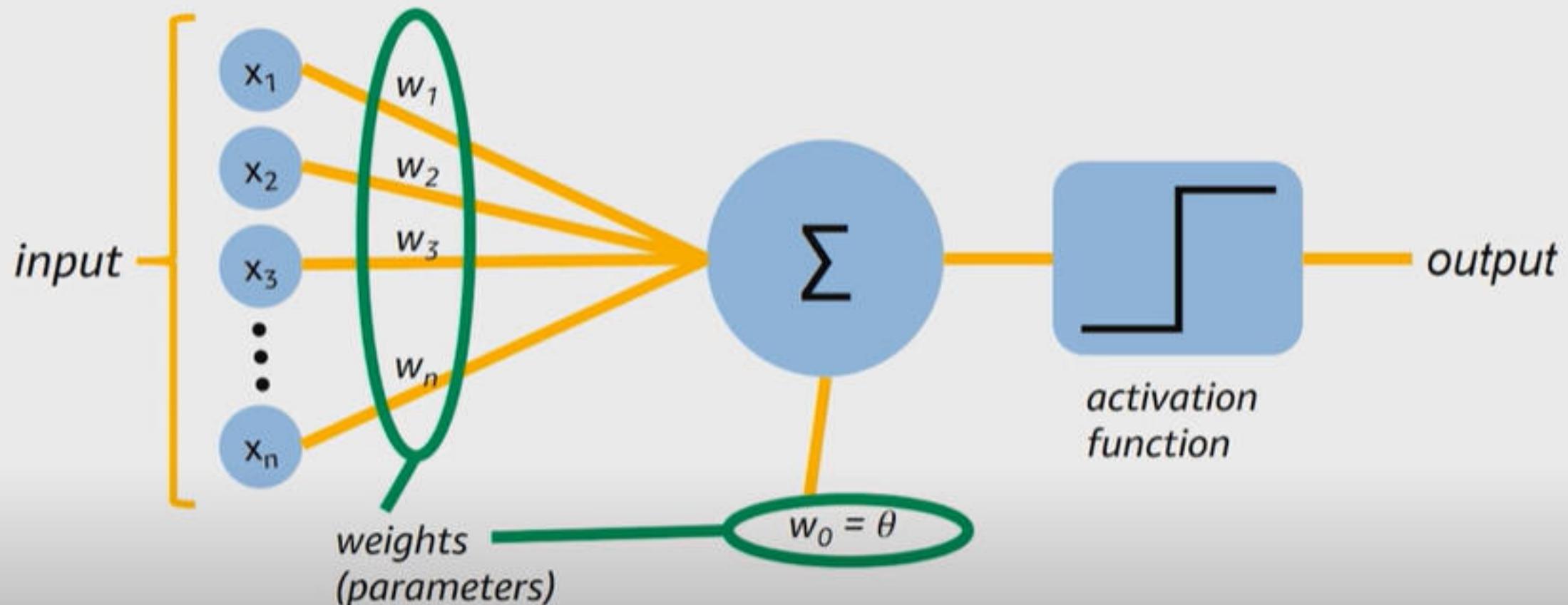


Kinesis Data  
Analytics

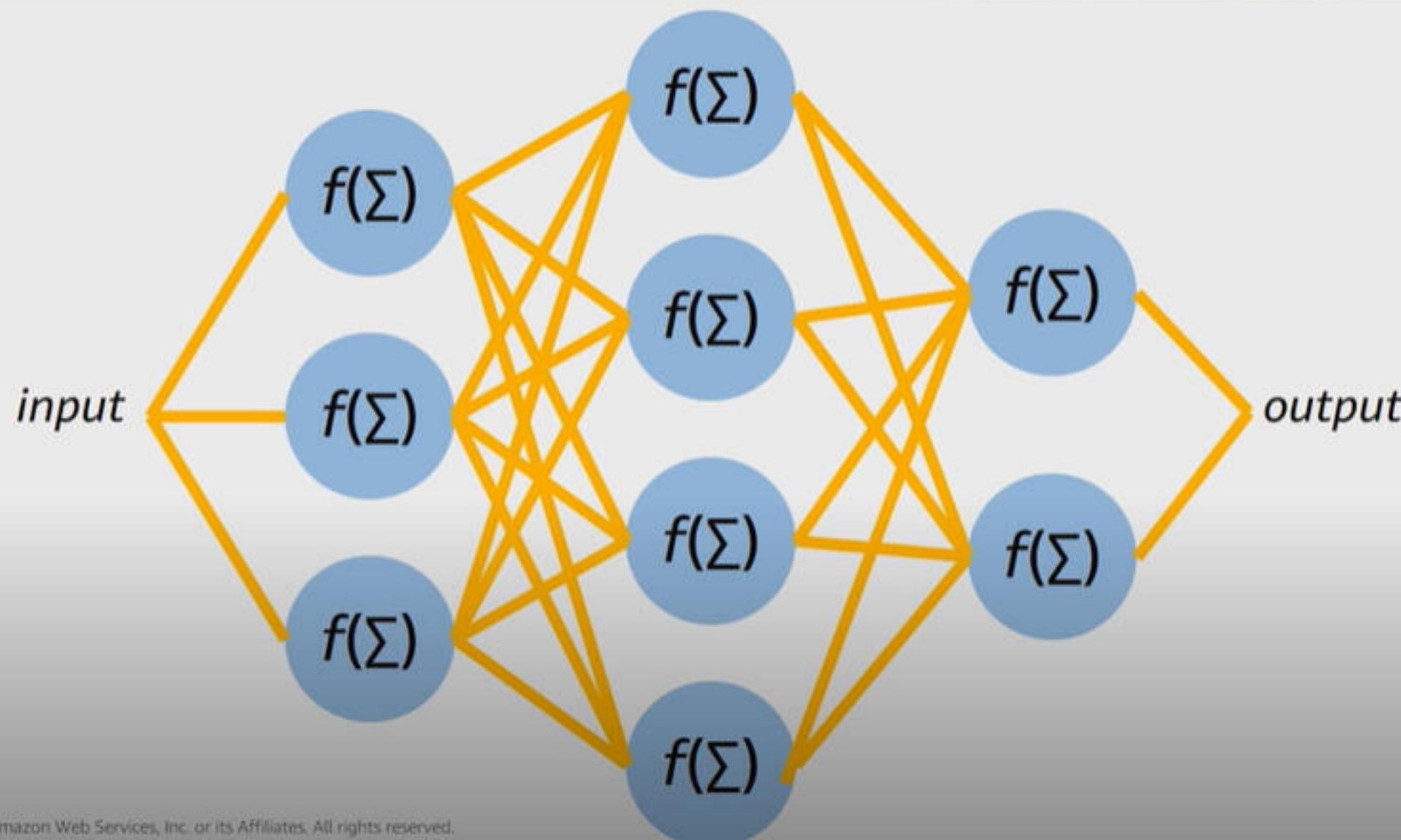
# Time Series Forecasting



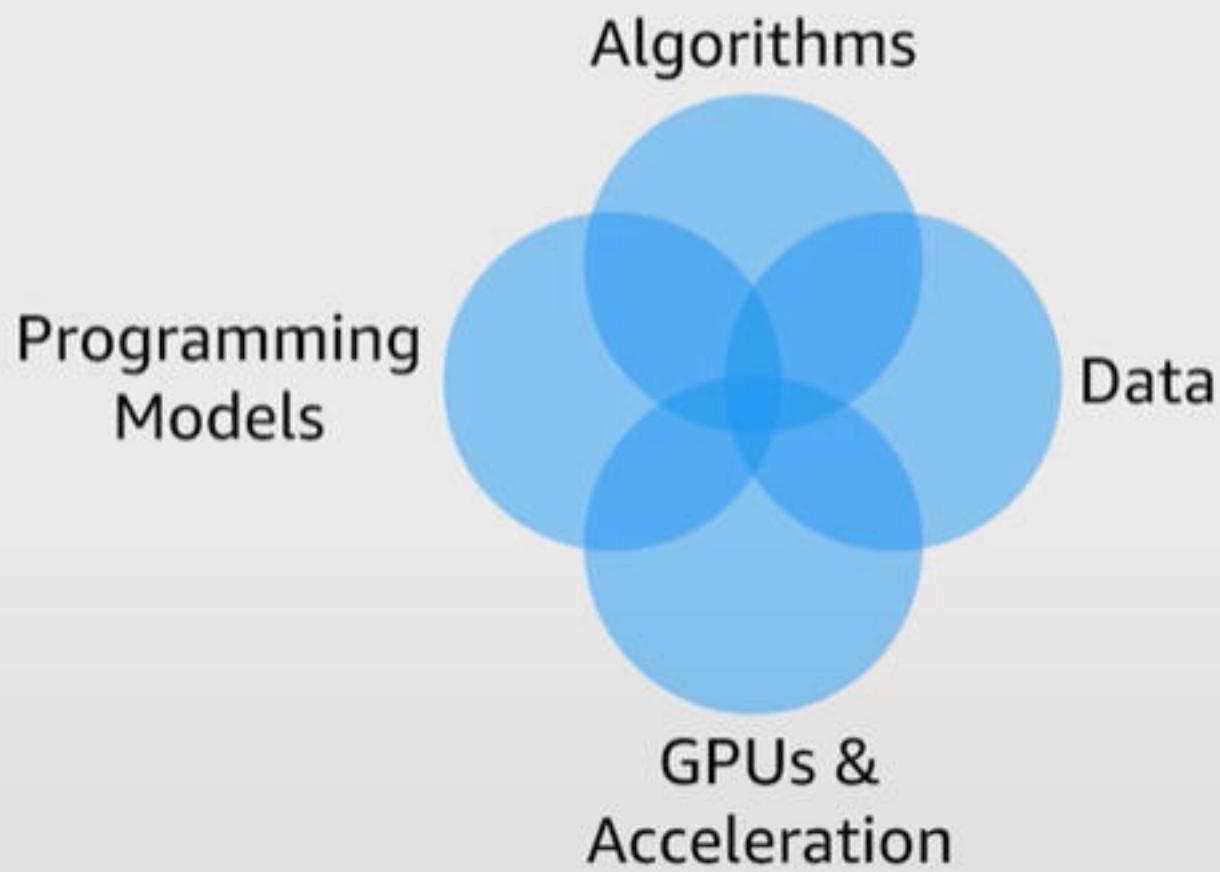
# Deep Learning: Neurons



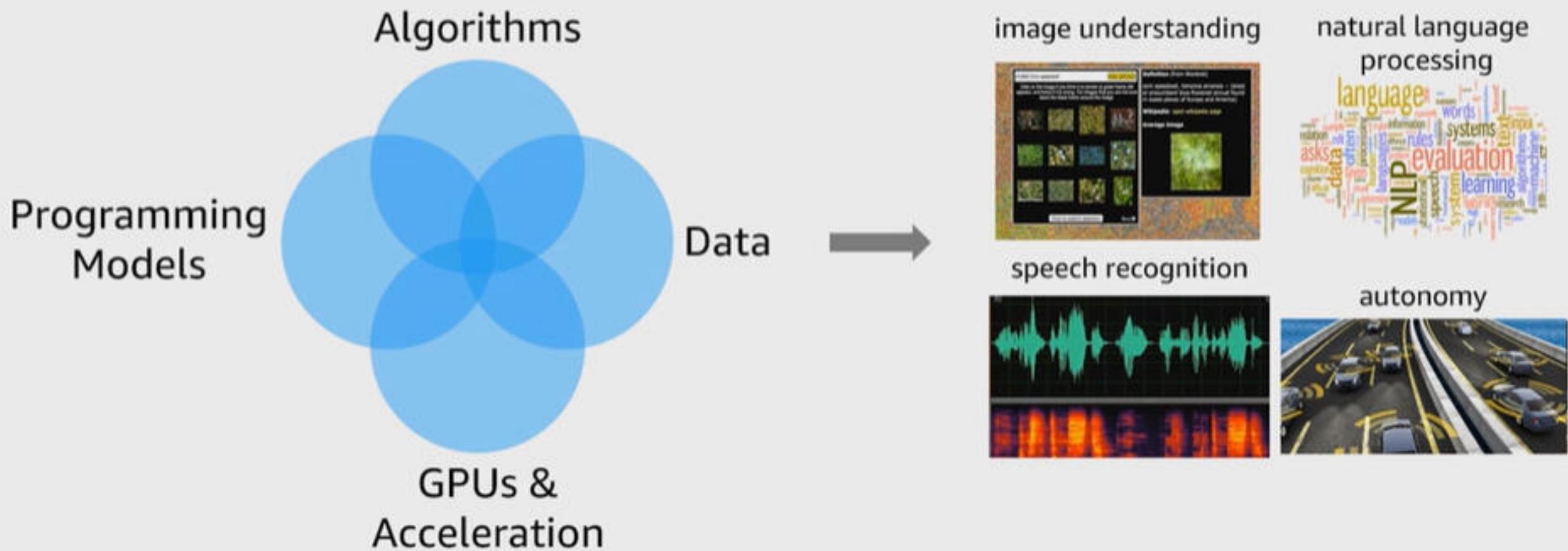
# Deep Learning: Neural Networks



# The Advent of Deep Learning

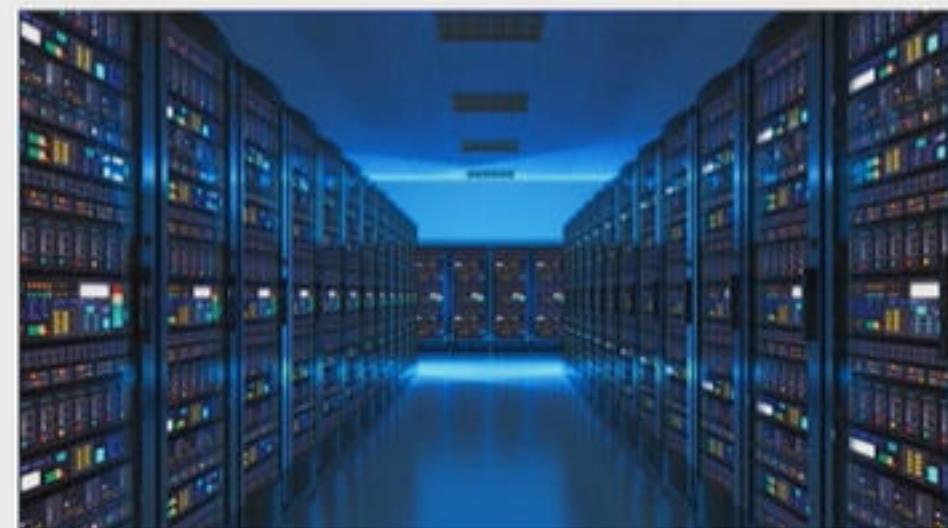
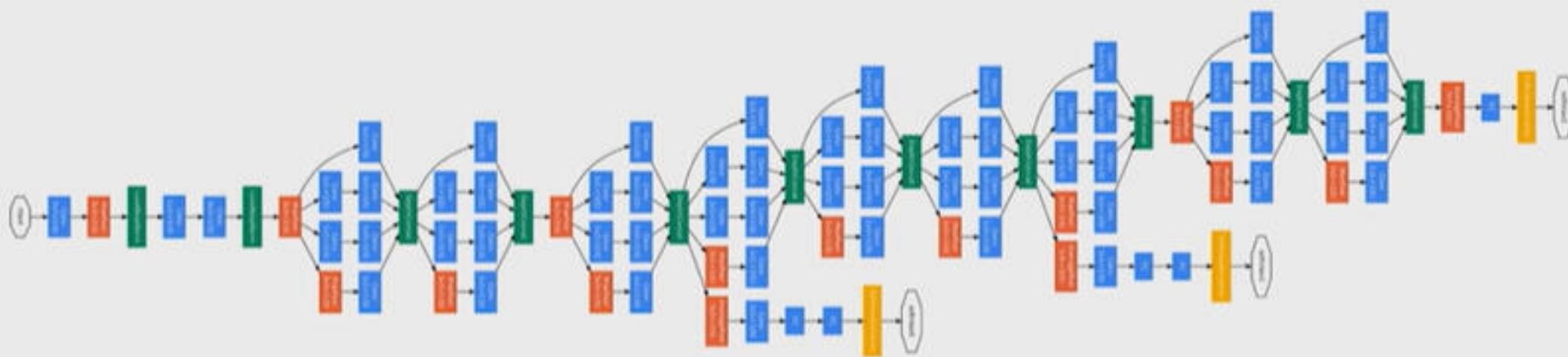


# The Advent of Deep Learning

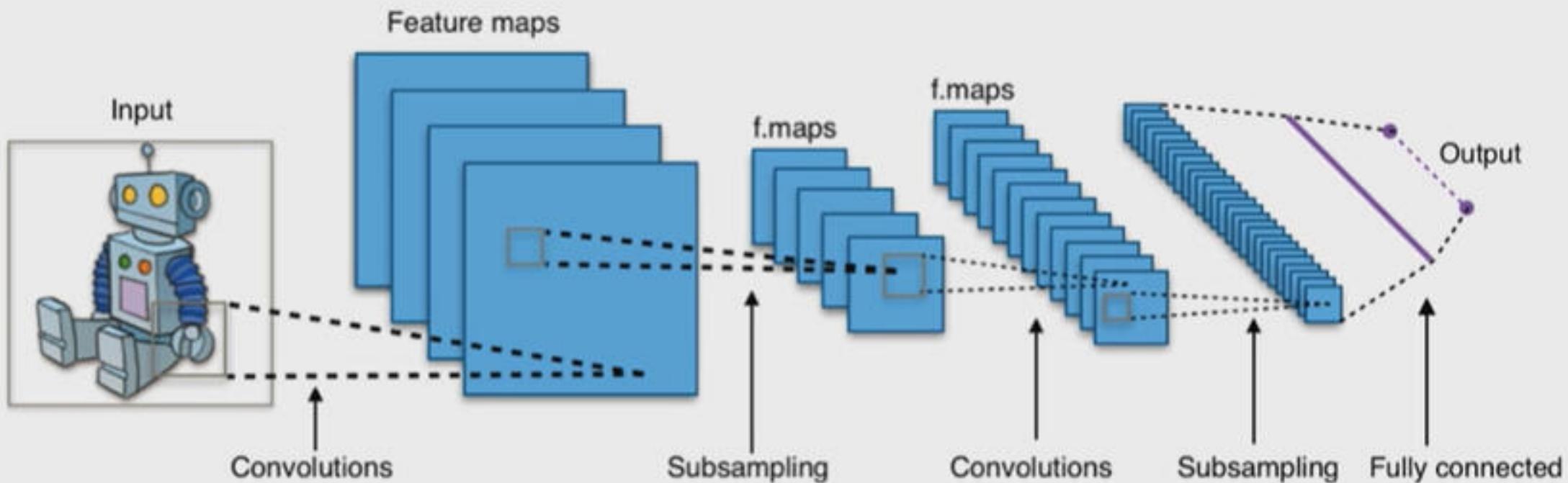


# Deep Neural Networks

aws training and certification



# Convolutional Neural Networks (CNNs)



# Convolutional Neural Networks (CNNs)



Object recognition, image classification



Semantic segmentation

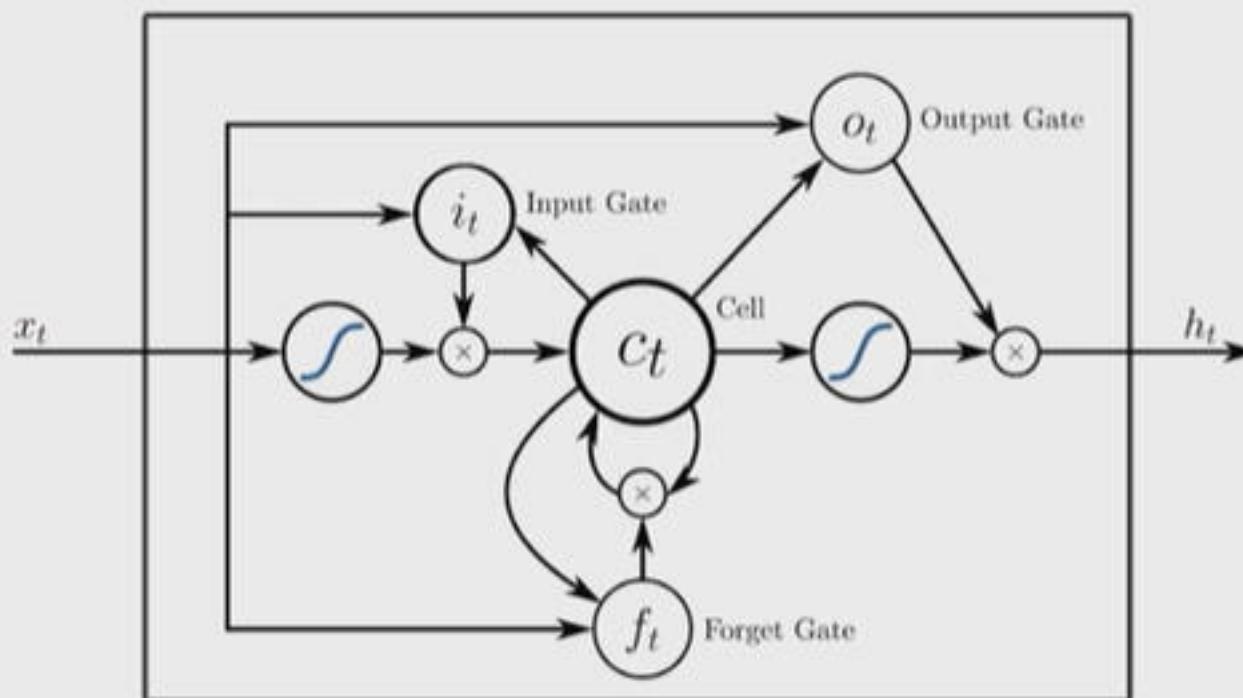
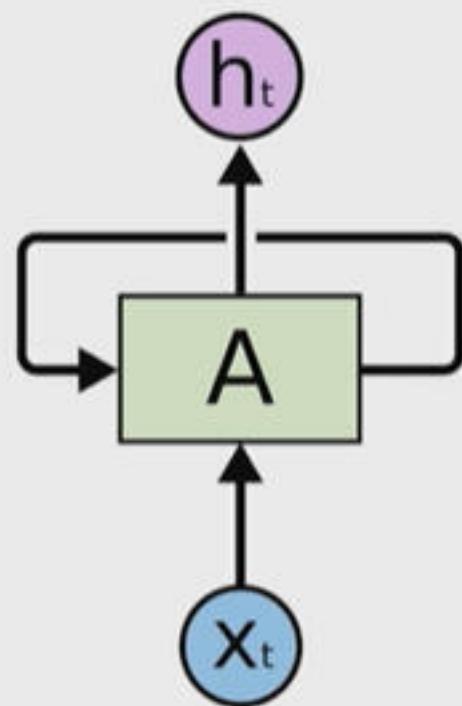


Artistic style transfer



Meow Generator

# Recurrent Neural Networks (RNNs)

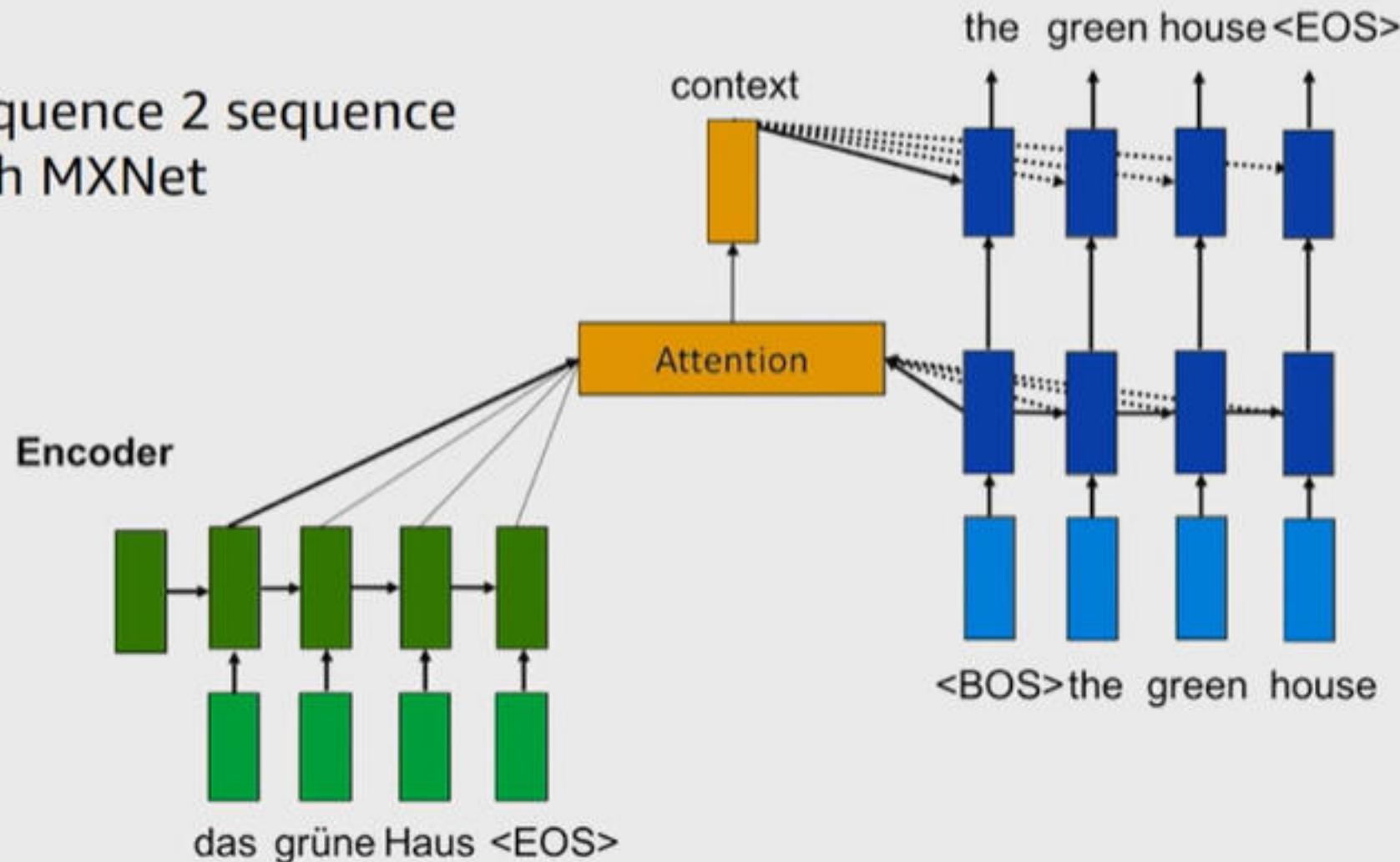


Long Short-Term Memory  
or LSTM network

# Recurrent Neural Networks (RNNs)



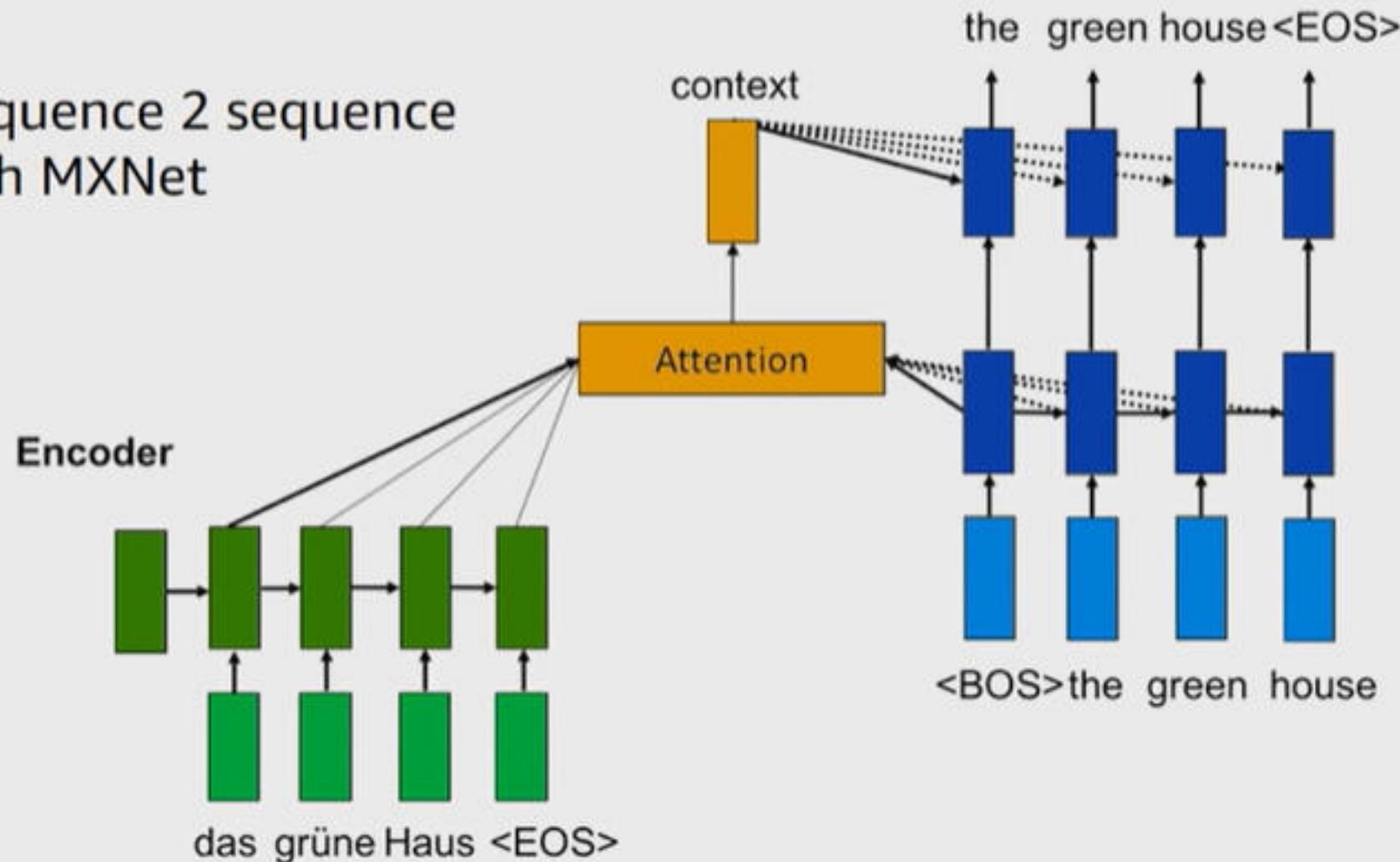
**Sockeye** – sequence 2 sequence  
modeling with MXNet



# Recurrent Neural Networks (RNNs)



**Sockeye** – sequence 2 sequence modeling with MXNet



# Different Network Topologies



A mostly complete chart of

## Neural Networks

- Backfed Input Cell
- Input Cell
- △ Noisy Input Cell
- Hidden Cell
- Probabilistic Hidden Cell
- △ Spiking Hidden Cell
- Output Cell
- Match Input Output Cell
- Recurrent Cell
- Memory Cell
- △ Different Memory Cell
- Kernel
- Convolution or Pool

©2016 Fjedor van Veen - asimovinstitute.org

Perceptron (P)



Feed Forward (FF)



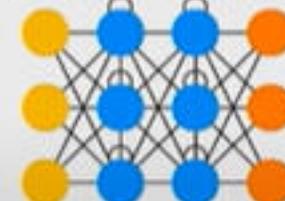
Radial Basis Network (RBF)



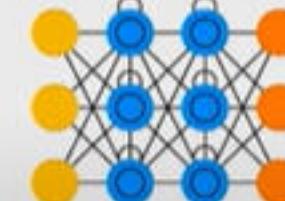
Deep Feed Forward (DFF)



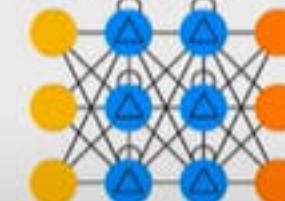
Recurrent Neural Network (RNN)



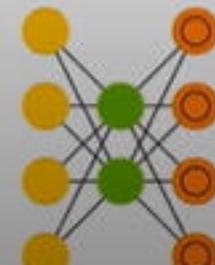
Long / Short Term Memory (LSTM)



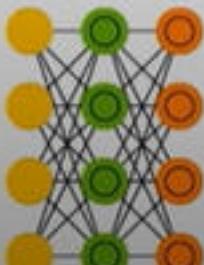
Gated Recurrent Unit (GRU)



Auto Encoder (AE)



Variational AE (VAE)



Denoising AE (DAE)



Sparse AE (SAE)



# Deep Learning Algorithms



## /Image Classification

- Convolutional Neural Network (ResNet)



## Sequence to Sequence (seq2seq)

- RNN for text summarization, translation, TTS



## Neural Topic Modelling (NTM)



## DeepAR Forecasting

- Time Series Prediction

