



THE
DEVELOPER'S
CONFERENCE

Building a Text Classifier with NaturalLanguage and CoreML

Alan Paiva

JULHO | 2019



 @ajeferson

ALAN PAIVA



iOS Software Engineer at ArcTouch,
Reactive Programming and Machine
Learning enthusiast, blog writer.
Power metal **headbanger** while doing
all of that.





We're hiring!

arctouch.com/careers/brazil

AGENDA

Intro

NaturalLanguage Framework

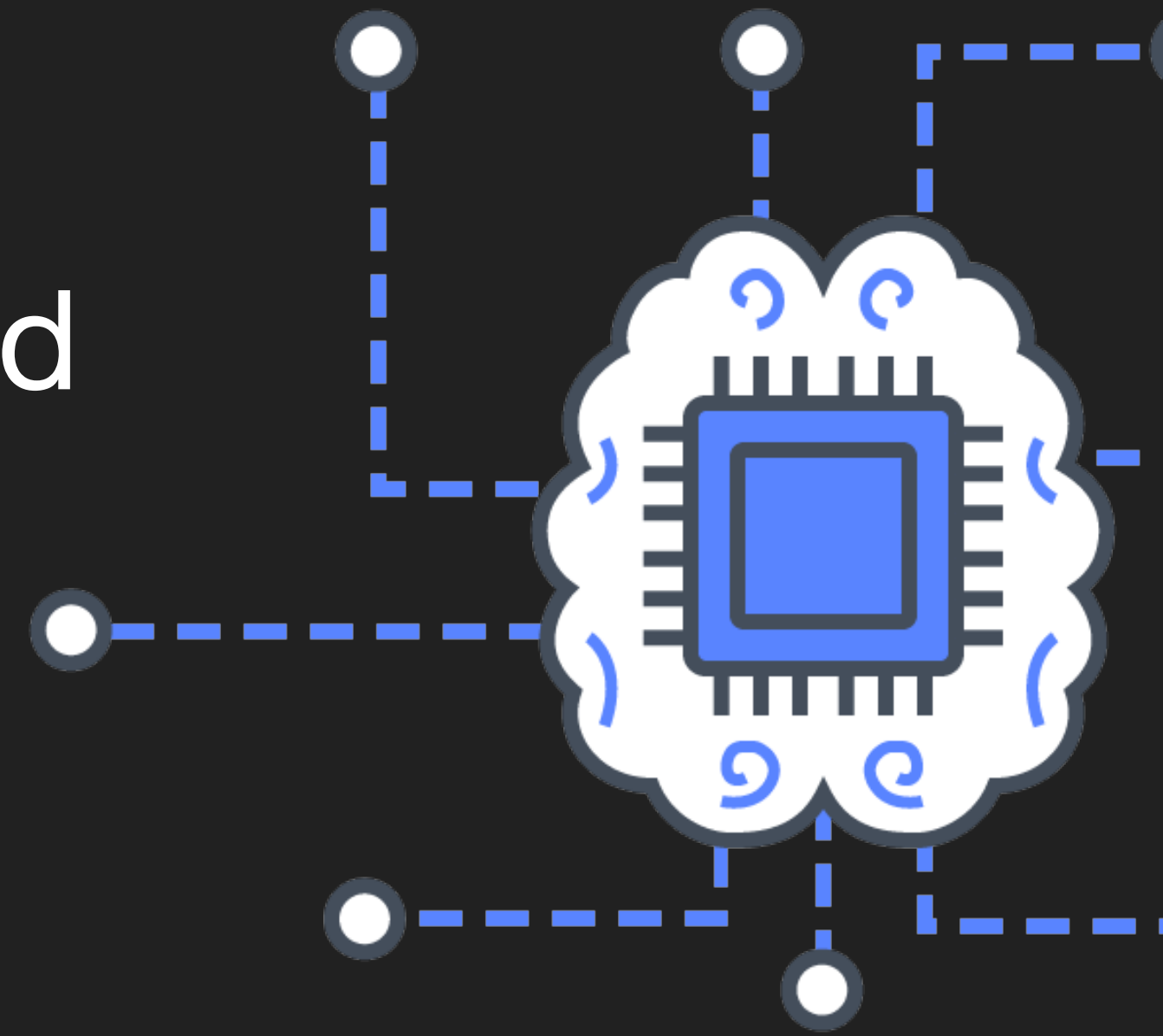
Machine Learning

Word Embeddings

Intro

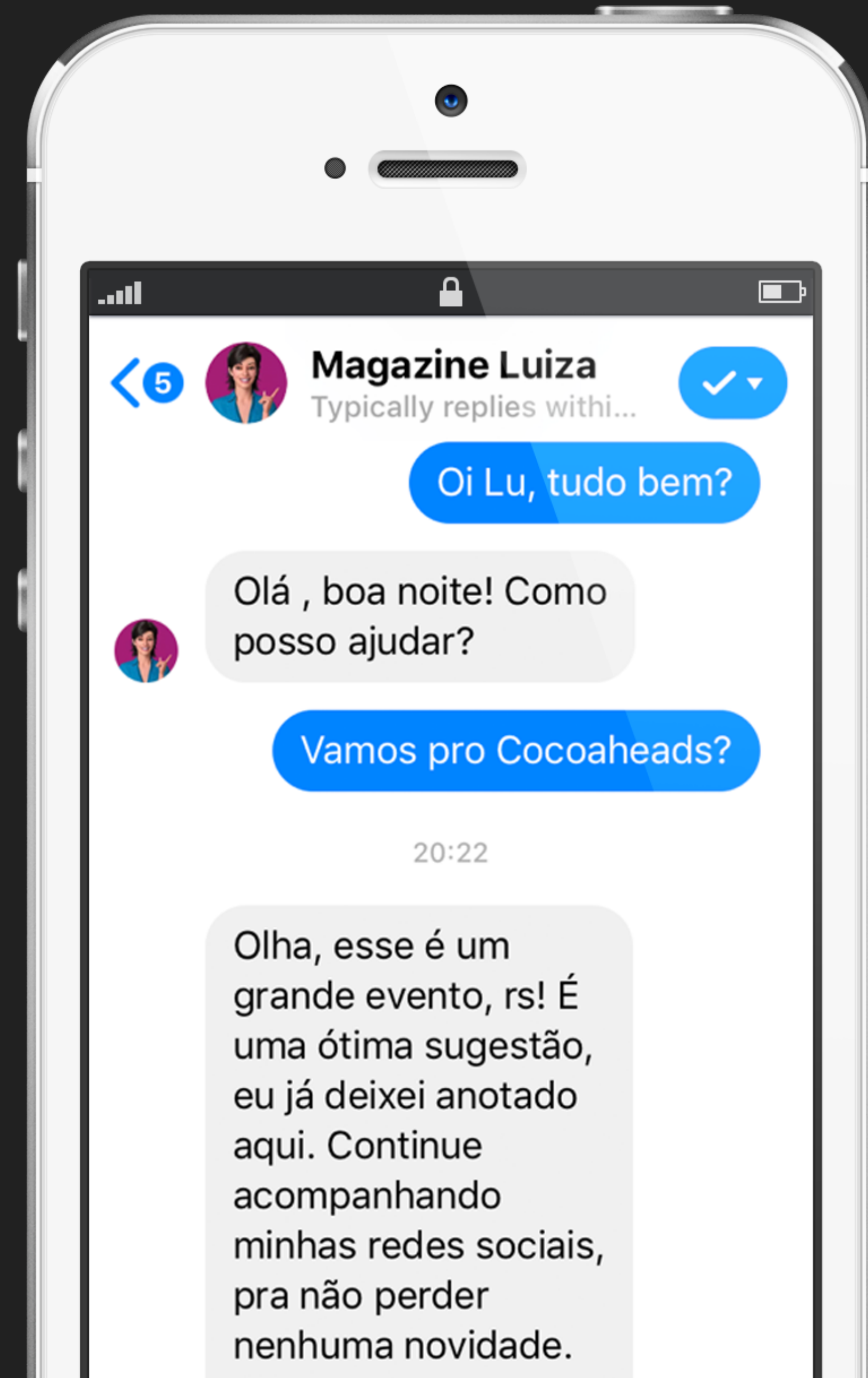
NLP

- In a wide sense, it refers to **any kind of natural language manipulation** performed by computers
- Field of **Machine Learning** that seeks for understanding human language



APPLICATIONS

- Text Categorization
- Machine Translation
- Chat Bots
- Virtual Assistants
- Sentiment Analysis
- SPAM Analysis



Natural Language

Framework

NATURAL LANGUAGE FRAMEWORK

Built-in Features

- Language Identification
- Tokenization
- POS-Tagging
- Named Entity Recognition

NATURAL LANGUAGE FRAMEWORK

- Integration with CreateML and CoreML
- Customizable Models



Text Classification

Word Tagging

Text Classification

TEXT

TOPIC

However, the site earlier announced that it was considering introducing restrictions on live-streaming in the wake of the Christchurch attacks. On Thursday, it also said that it would ban white nationalism and separatism from the site.

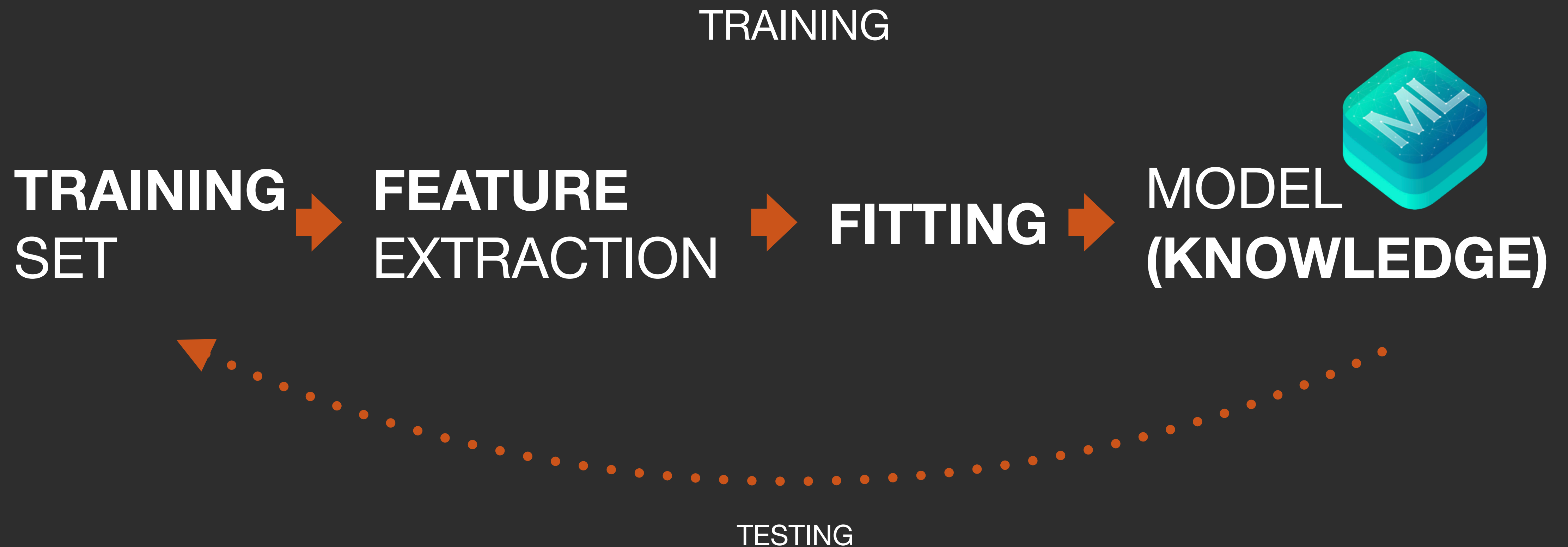
NEWS

And God said, "Let there be light," and there was light. 4 God saw that the light was good, and he separated the light from the darkness. 5 God called the light "day," and the darkness he called "night." And there was evening, and there was morning—the first day.

RELIGION

Machine Learning

Supervised Learning



Supervised Learning

PREDICTING



Supervised Learning

NaturalLanguage Framework

TRAINING

TRAINING DATA



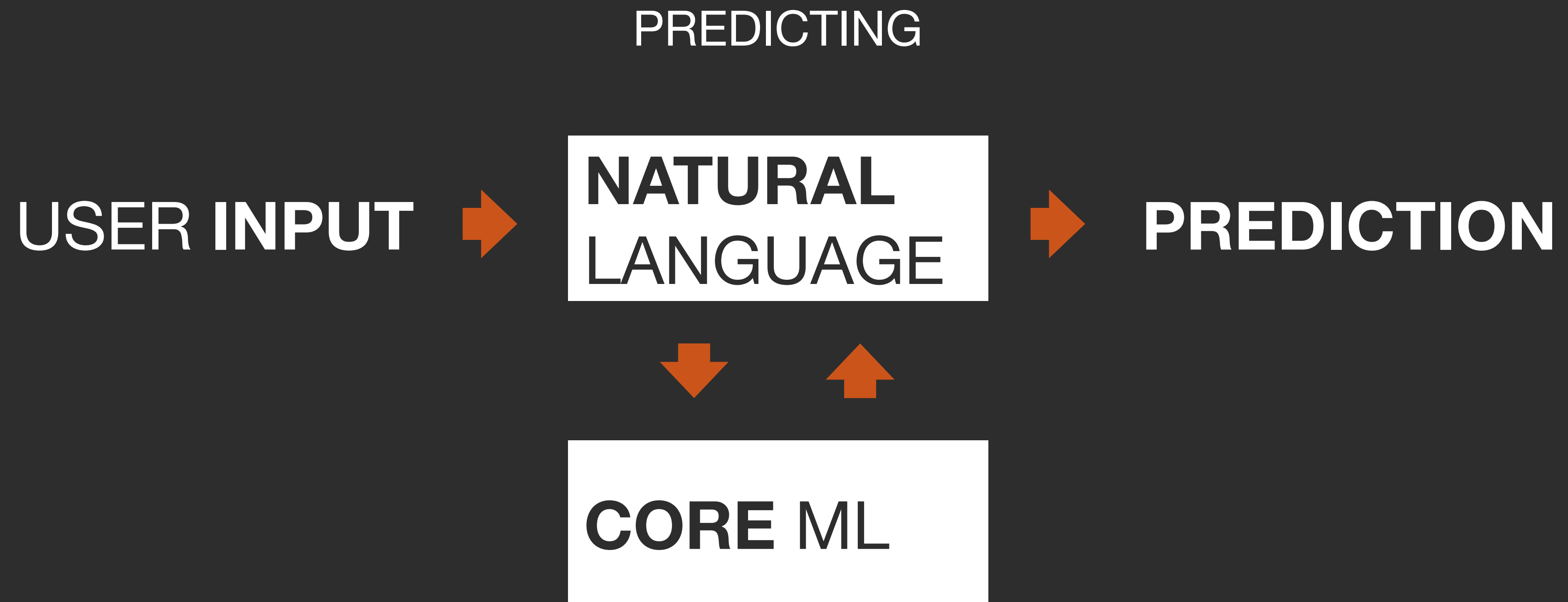
CREATE ML



NATURAL
LANGUAGE

Supervised Learning

NaturalLanguage Framework



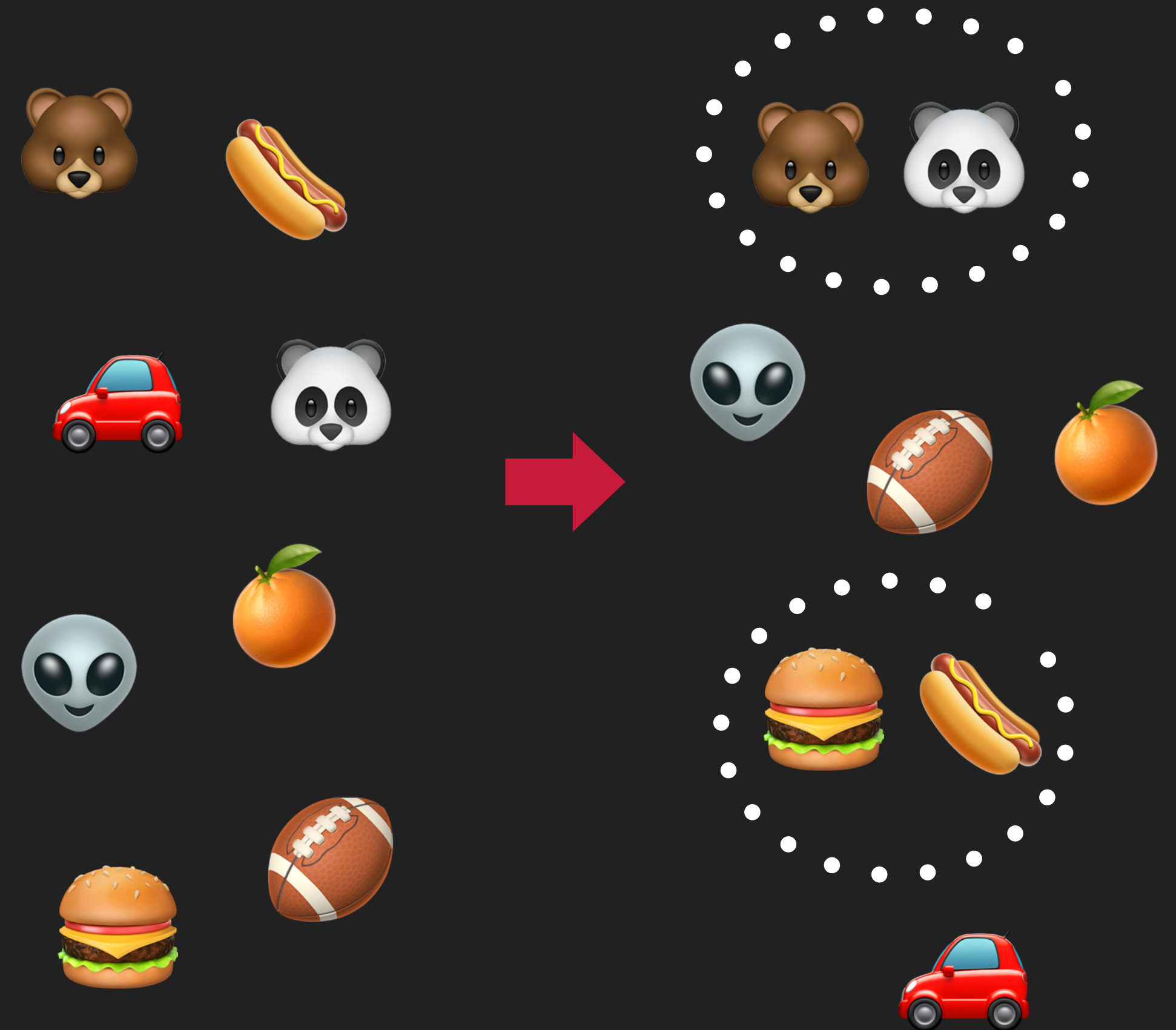
DEMO



Word Embeddings

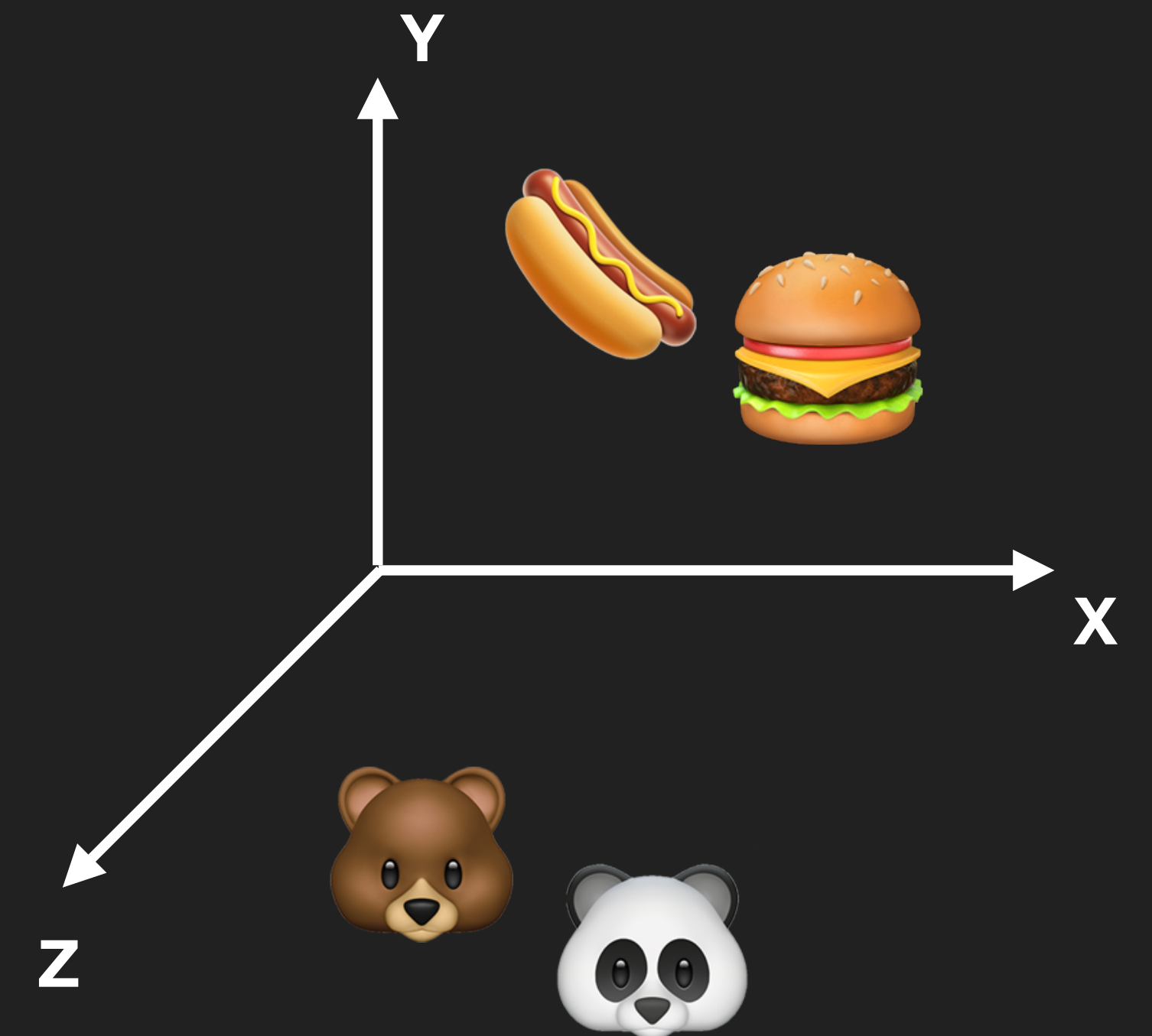
WORD EMBEDDINGS

- Groups words that are **similar in meaning**
- Map strings to **vectors** in a n-dimensional space



WORD EMBEDDINGS

	9.874	6.409	7.329
	9.539	6.087	6.997
...
	-2.341	-6.455	-1.532
	-2.312	-6.502	-1.704



TRANSFER LEARNING

- Usage of WordEmbeddings to **improve models**
- **Previous knowledge** of the language (pre-trained model)
- **Less training data** required
- Types: **static, dynamic and custom**

TRANSFER LEARNING

```
// Creating the model

let params = MLTextClassifier.ModelParameters(algorithm: .maxEnt(revision: 1))

let model = try MLTextClassifier(trainingData: trainingSet,
                                textColumn: "text",
                                labelColumn: "label",
                                parameters: params)
```

TRANSFER LEARNING

```
// Creating the model

let params = MLTextClassifier.ModelParameters(algorithm: .transferLearning(.staticEmbedding,
                                                                    revision: 1))

let model = try MLTextClassifier(trainingData: trainingSet,
                                textColumn: "text",
                                labelColumn: "label",
                                parameters: params)
```


DEMO



TRANSFER LEARNING

Project: TopicClassifier

Model Sources:

- TopicClassifier 1 - MaxEnt
- TopicClassifier 2 - StaticEmbeddings
- TopicClassifier 3 - DynamicEmbeddings

Input: 4 Classes

Accuracy: 100% Training, 74% Validation, 76% Testing

Output: 76 KB

▼ Data Inputs

Data Type	Count	Unit	Source
Training Data	400	Items	training
Validation Data	Auto		Automatic
Testing Data	100	Items	test

▼ Parameters

Algorithm:

- Maximum entropy
- Conditional random field
- Transfer learning

TRANSFER LEARNING

The screenshot displays a machine learning interface for a TopicClassifier. The main dashboard shows the following metrics:

Input	Accuracy	Output
4 Classes	95% Training	1,1 MB
	87% Validation	
	84% Testing	

The interface also shows the following data inputs and parameters:

Data Inputs

Training Data	Validation Data	Testing Data
400 Items	Auto	100 Items
training	Automatic	test

Parameters

- Algorithm: Maximum entropy, Conditional random field, Transfer learning
- Feature Extractor: Static embedding

TRANSFER LEARNING

Project

- TopicClassifier

Model Sources

- TopicClassifier 1 - MaxEnt
- TopicClassifier 2 - StaticEmbeddings
- TopicClassifier 3 - DynamicEmbeddings**

Input

4 Classes

Accuracy

82% Training

90% Validation

73% Testing

Output

1,5 MB

▼ Data Inputs

Training Data

400 Items

training

Validation Data

Auto

Automatic

Testing Data

100 Items

test

▼ Parameters

Algorithm

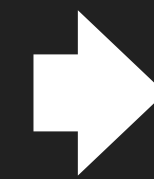
- Maximum entropy
- Conditional random field
- Transfer learning**

Feature Extractor

Dynamic embedding

RECAP

Natural Language
Machine Learning

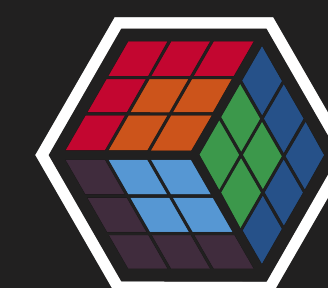


Text classifier

- WordEmbeddings
- Transfer learning

Thank
you

GITHUB REPO



THE
DEVELOPER'S
CONFERENCE