

TRAINING DEEP LEARNING MODELS AT SCALE USING KUBERNETES

Mitul Tiwari and Deepak Bobbarjung

 Passage AI

Introductions



Mitul Tiwari

CTO and Cofounder at Passage AI
San Francisco Bay Area | Computer Software

Current Passage AI, Forbes Technology Council
Previous LinkedIn, Kosmix, Google
Education The University of Texas at Austin



Deepak Bobbarjung

Founding Engineer at Passage AI
San Francisco Bay Area
| Information Technology and Services

Current Passage AI
Previous EMC, Maginatics, Inc (Acquired by EMC), VMware
Education Purdue University

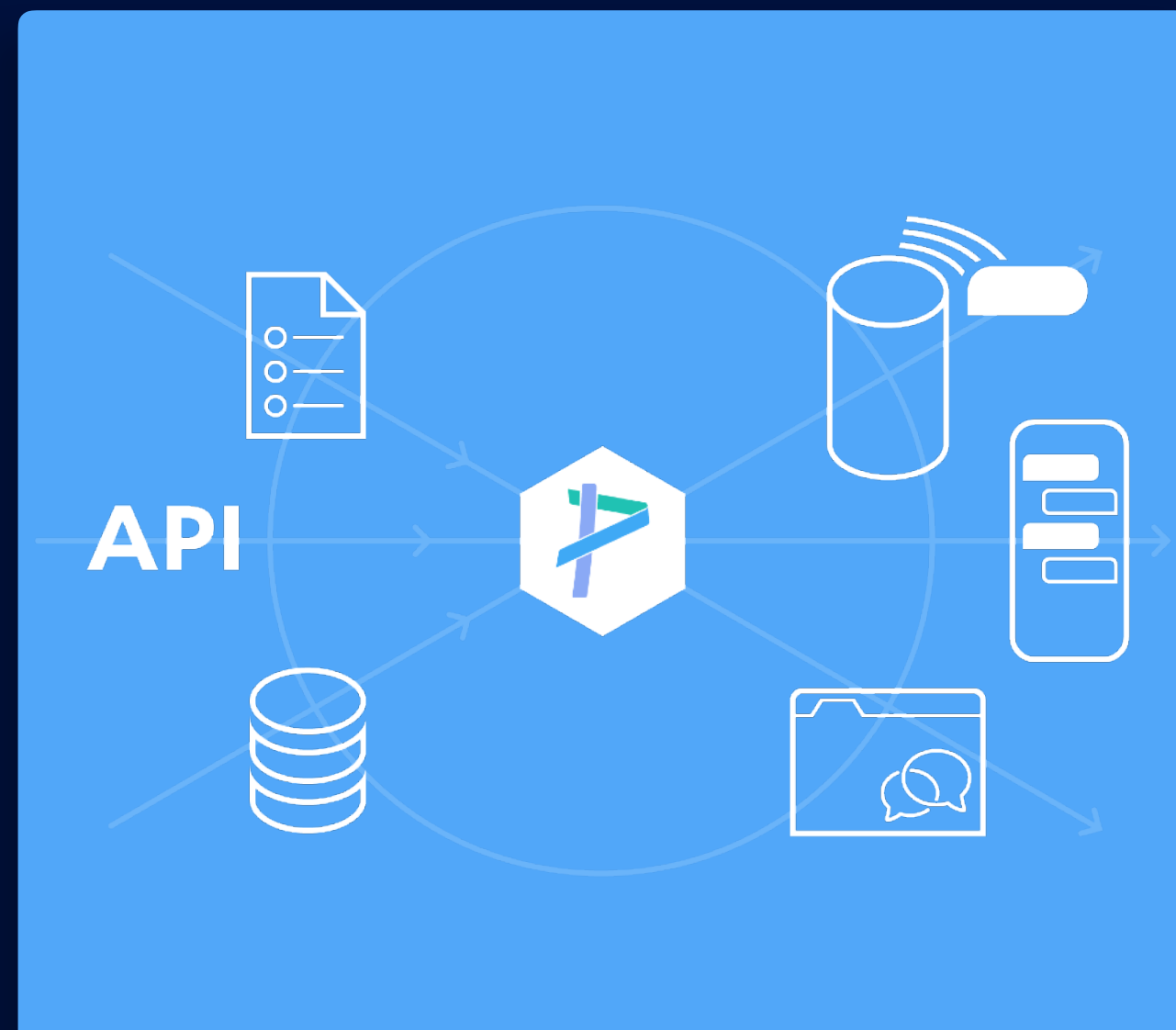
Outline

Conversational AI and Deep Learning

Need for a Jobs framework on Kubernetes

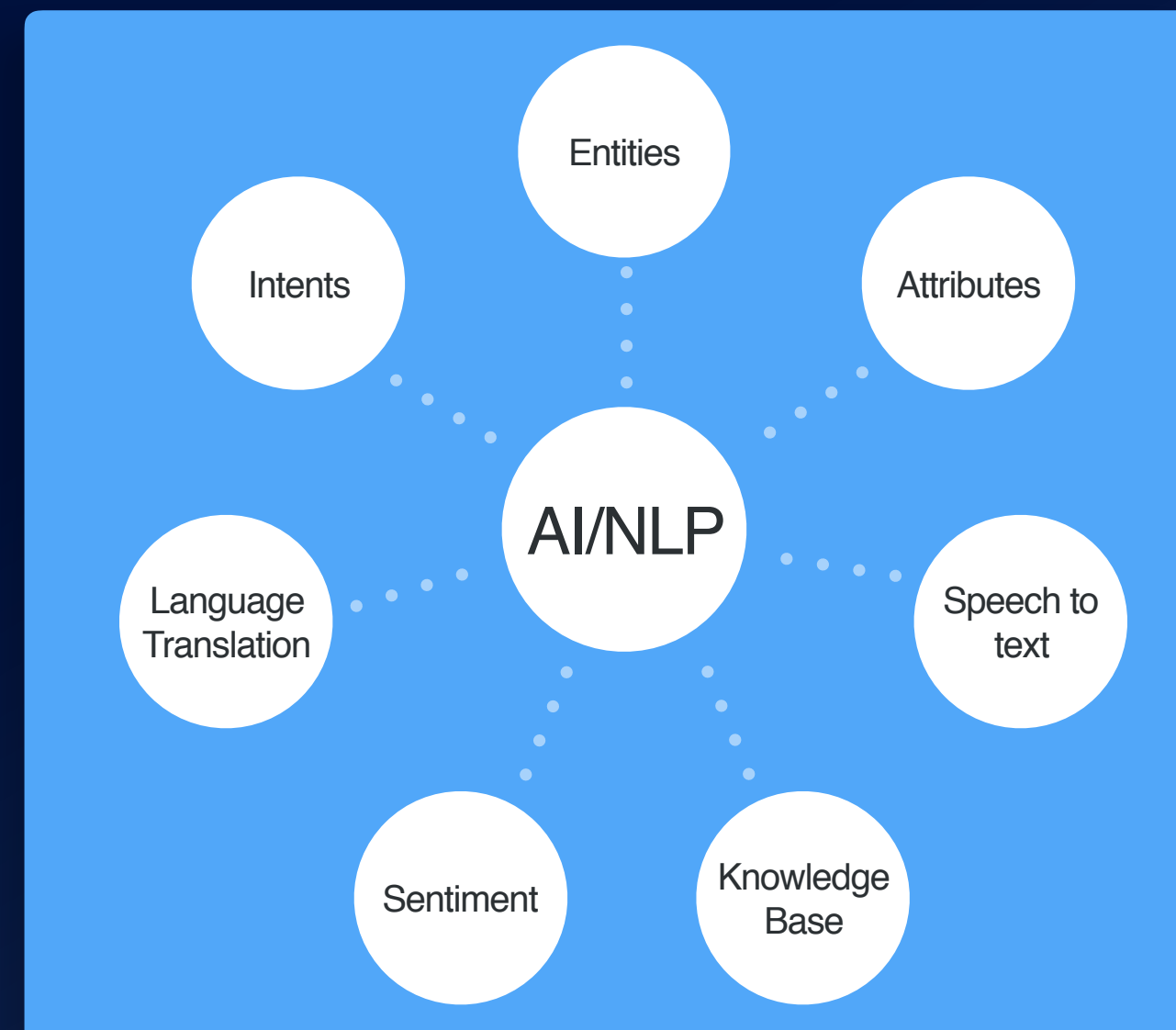
Our Jobs architecture

Our Conversational AI Platform



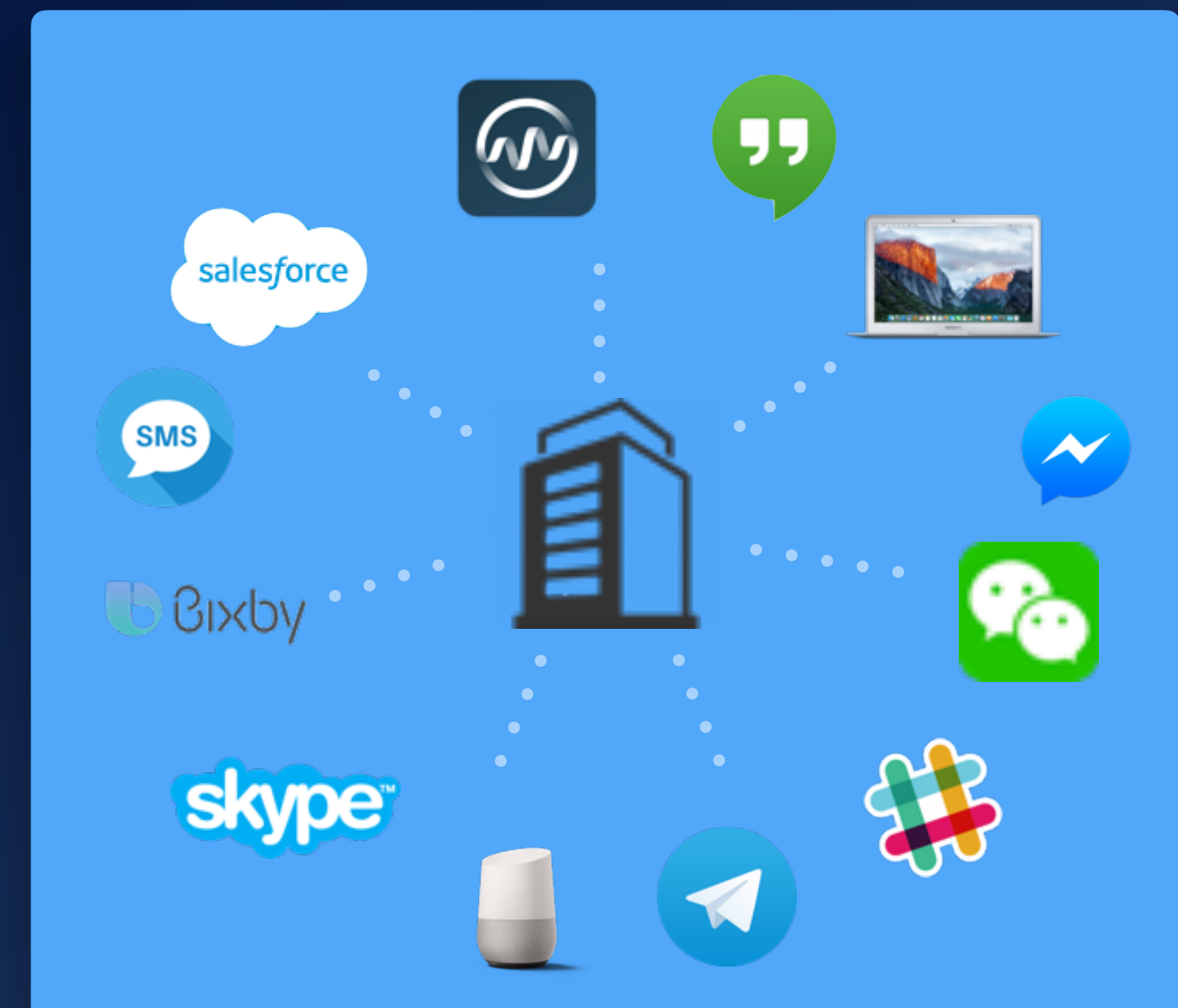
The diagram illustrates the Bot Builder process. A central hexagonal icon with a stylized 'P' and 'A' is connected by arrows to various components: a document icon labeled 'API', a database icon, a folder icon, a smartphone icon, and a server rack icon. The entire process is set against a light blue background.

Bot Builder
No Coding Required.



The diagram illustrates the Bot Training process. A central circle labeled 'AI/NLP' is connected by dotted lines to seven surrounding circles: 'Intents', 'Entities', 'Attributes', 'Speech to text', 'Knowledge Base', 'Sentiment', and 'Language Translation'. The entire process is set against a light blue background.

Bot Training
#1 AI/NLP Model.

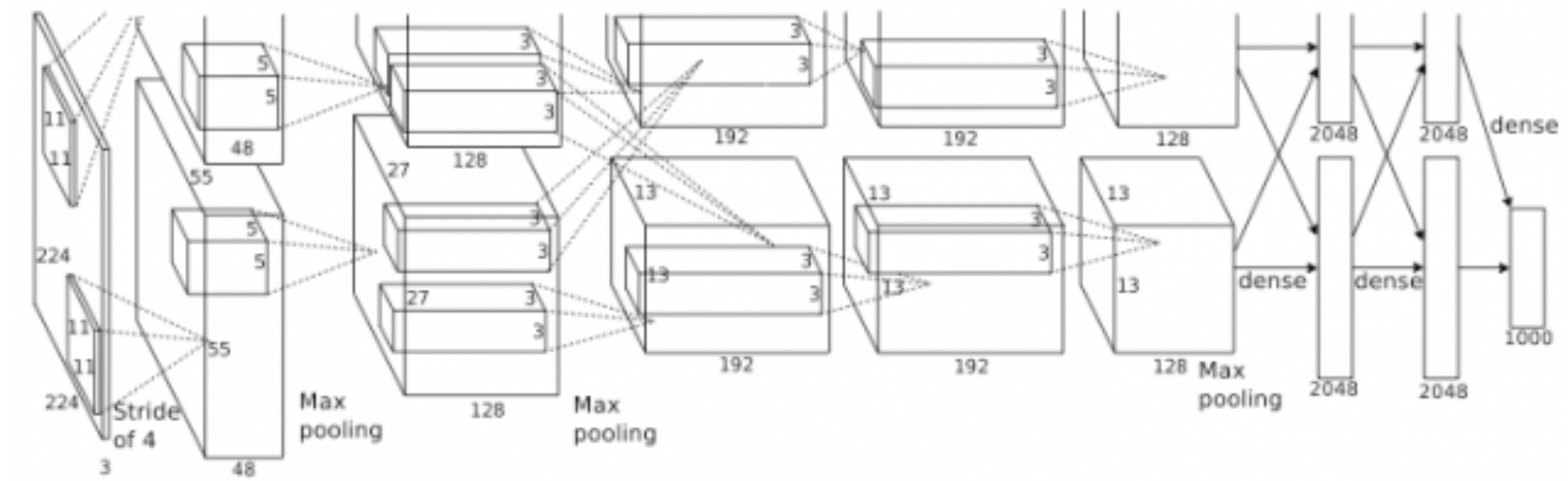
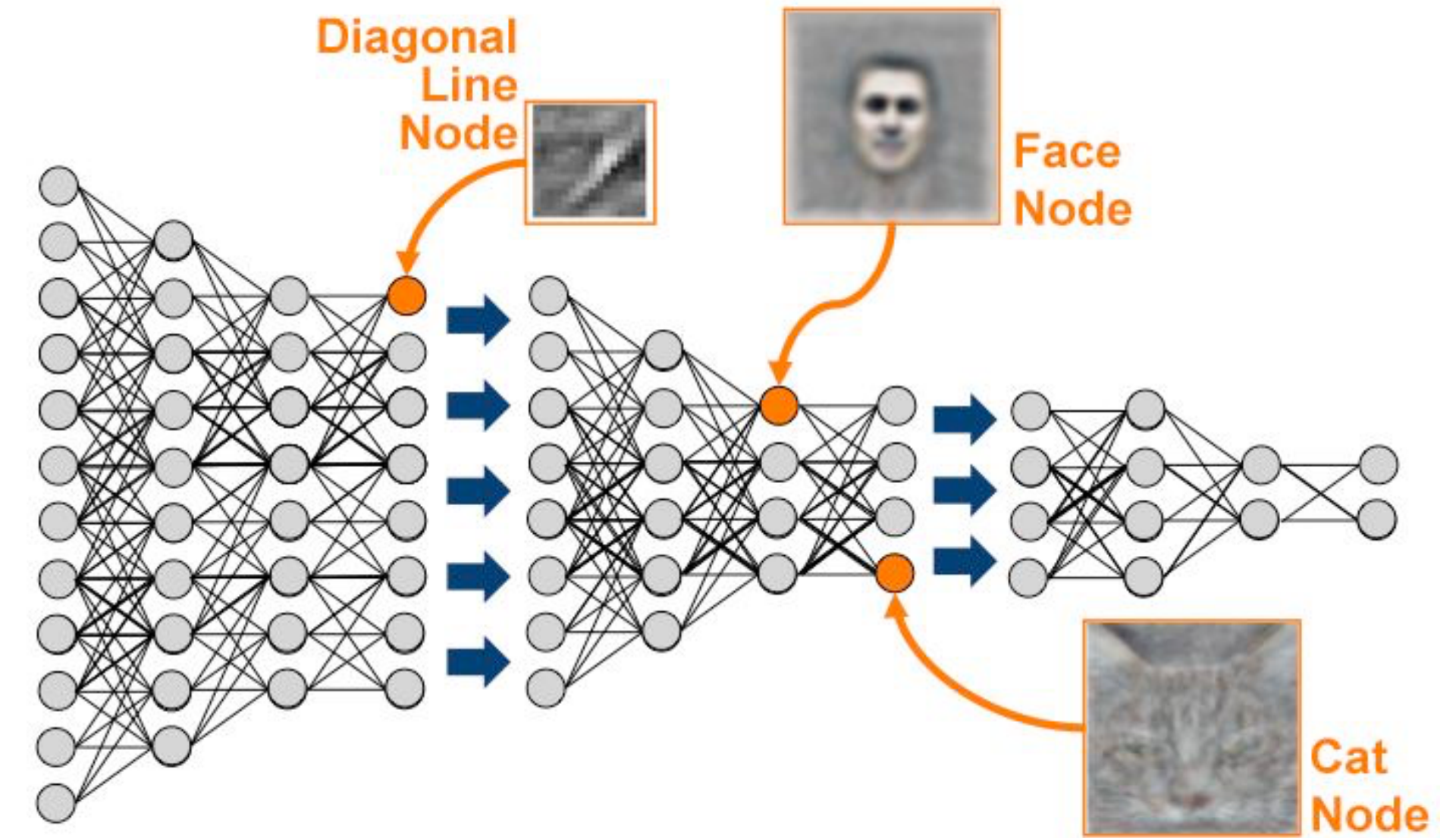


The diagram illustrates the Bot Deployment process. A central building icon is connected by dotted lines to various communication channels and devices: 'salesforce', 'SMS', 'Bixby', 'skype', a smart speaker, a laptop, a speech bubble, a chat bubble, a Telegram icon, and a colorful puzzle piece icon. The entire process is set against a light blue background.

Bot Deployment
Build Once, Deploy Everywhere.

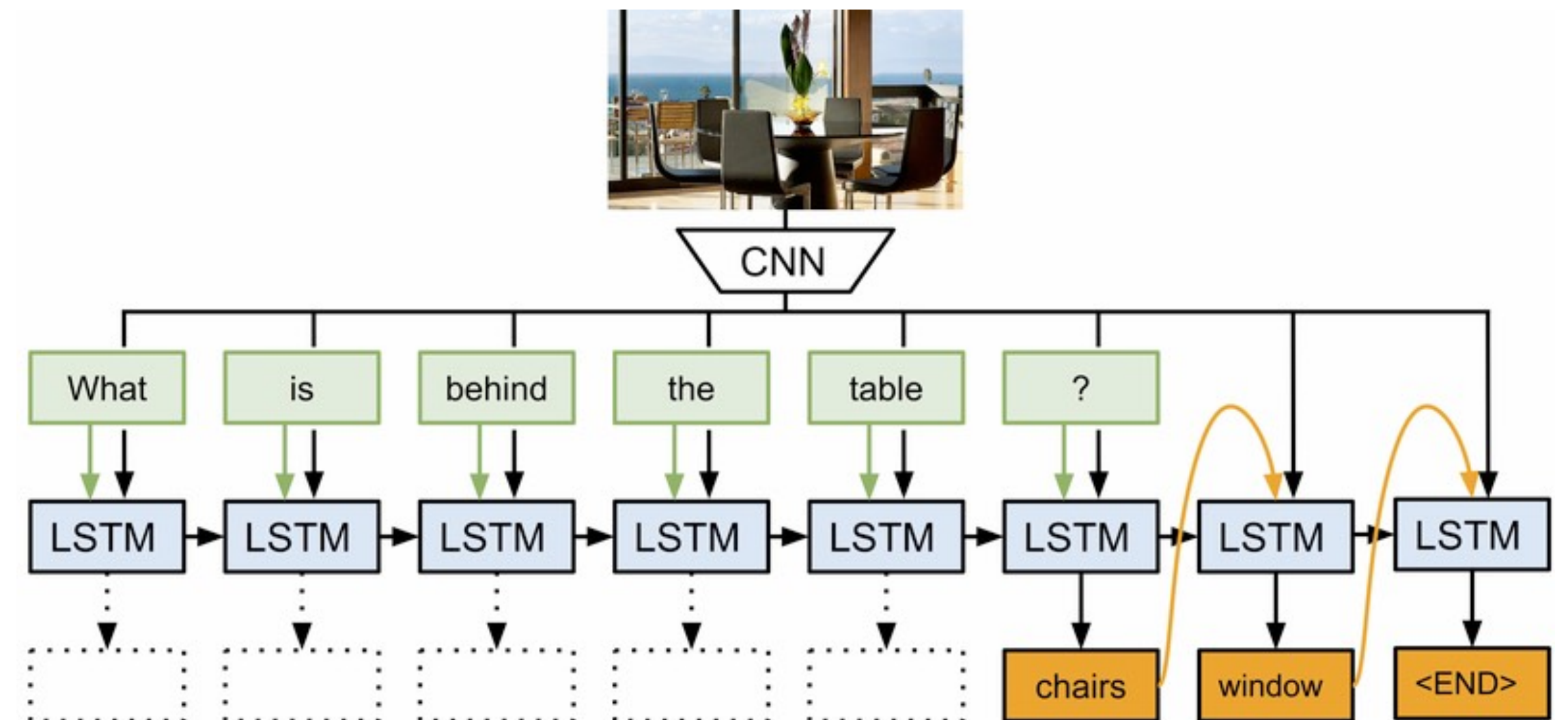
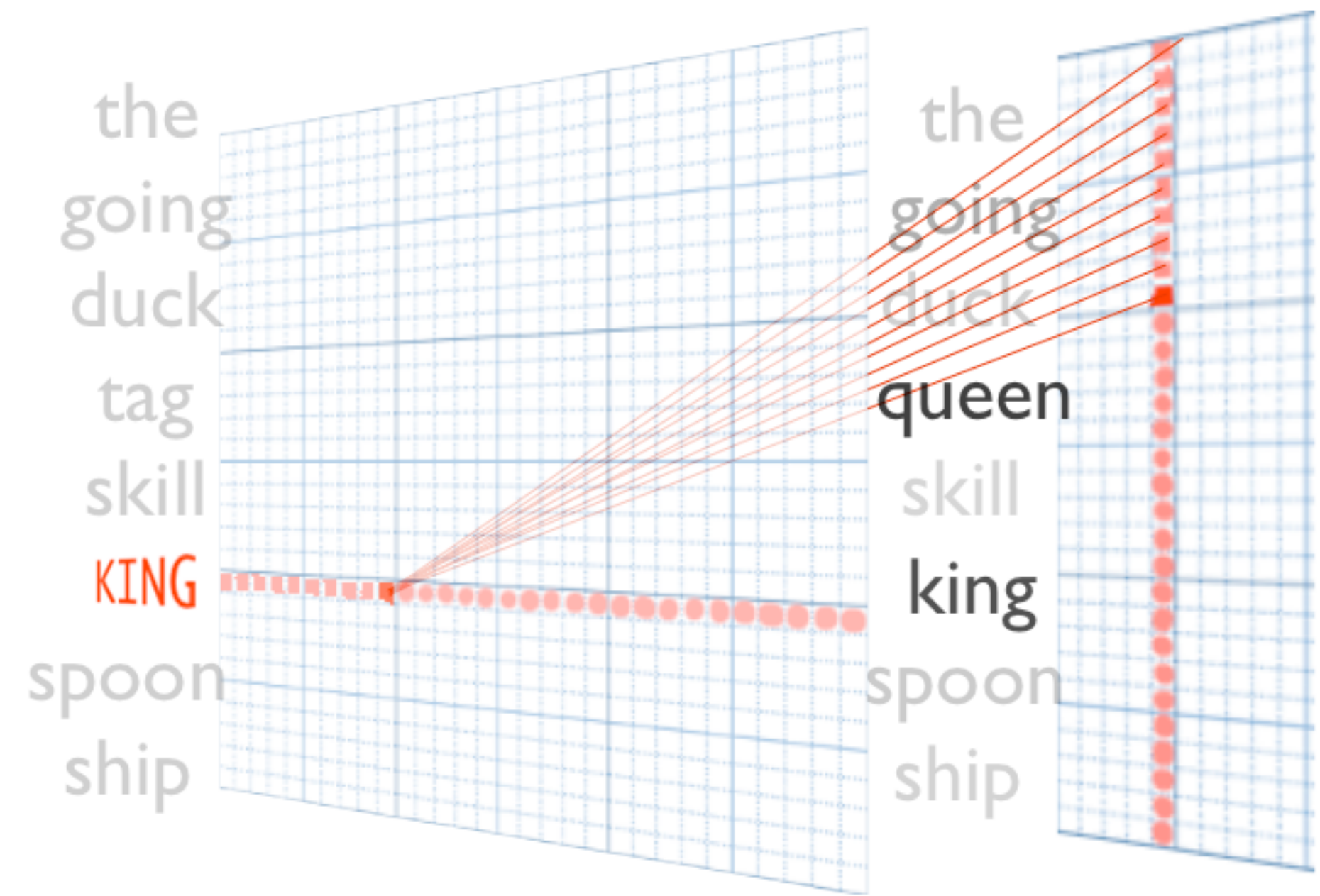
Deep Learning

- Traditional Machine Learning
 - Human designed features and representations
 - Optimize weights to combine
- Deep Learning
 - Deep Neural Network
 - Learn good features and multiple levels of representations



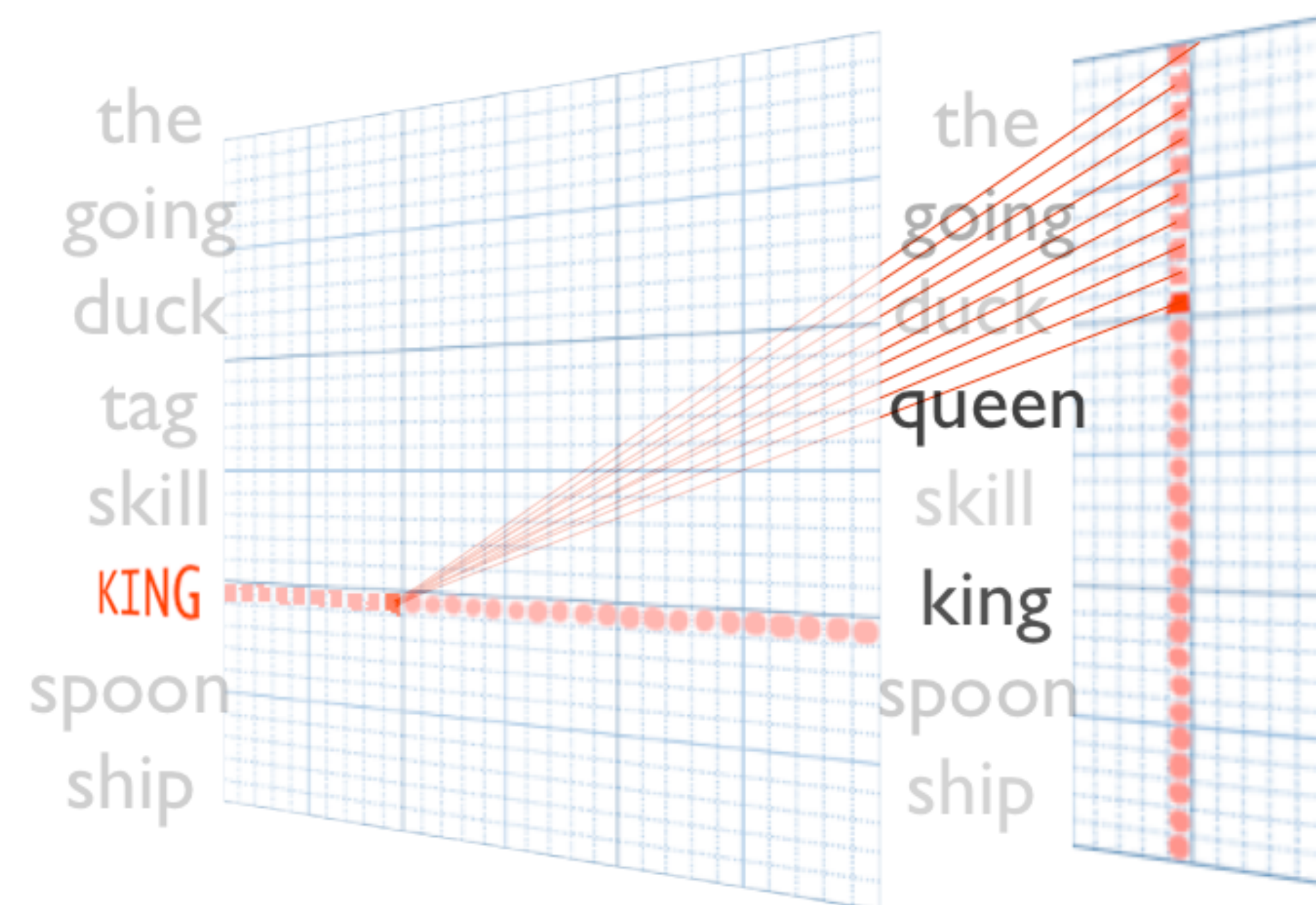
Deep Learning for NLP

- Language Translation
- Image Captioning
- Text Summarization
- Parts-of-speech Tagging
- Named Entity Recognition
- Natural Language Generation
- Question-Answering
- Optical Character Recognition
- Speech Recognition
- Machine Reading Comprehension



Neural Network for Word Embedding

- Word Embedding: Word2Vec
- Embed words in continuous vector space
- Semantically similar words are mapped to nearby points
- Enables powerful operations
- “King” - “Man” + “Woman” -> “Queen”



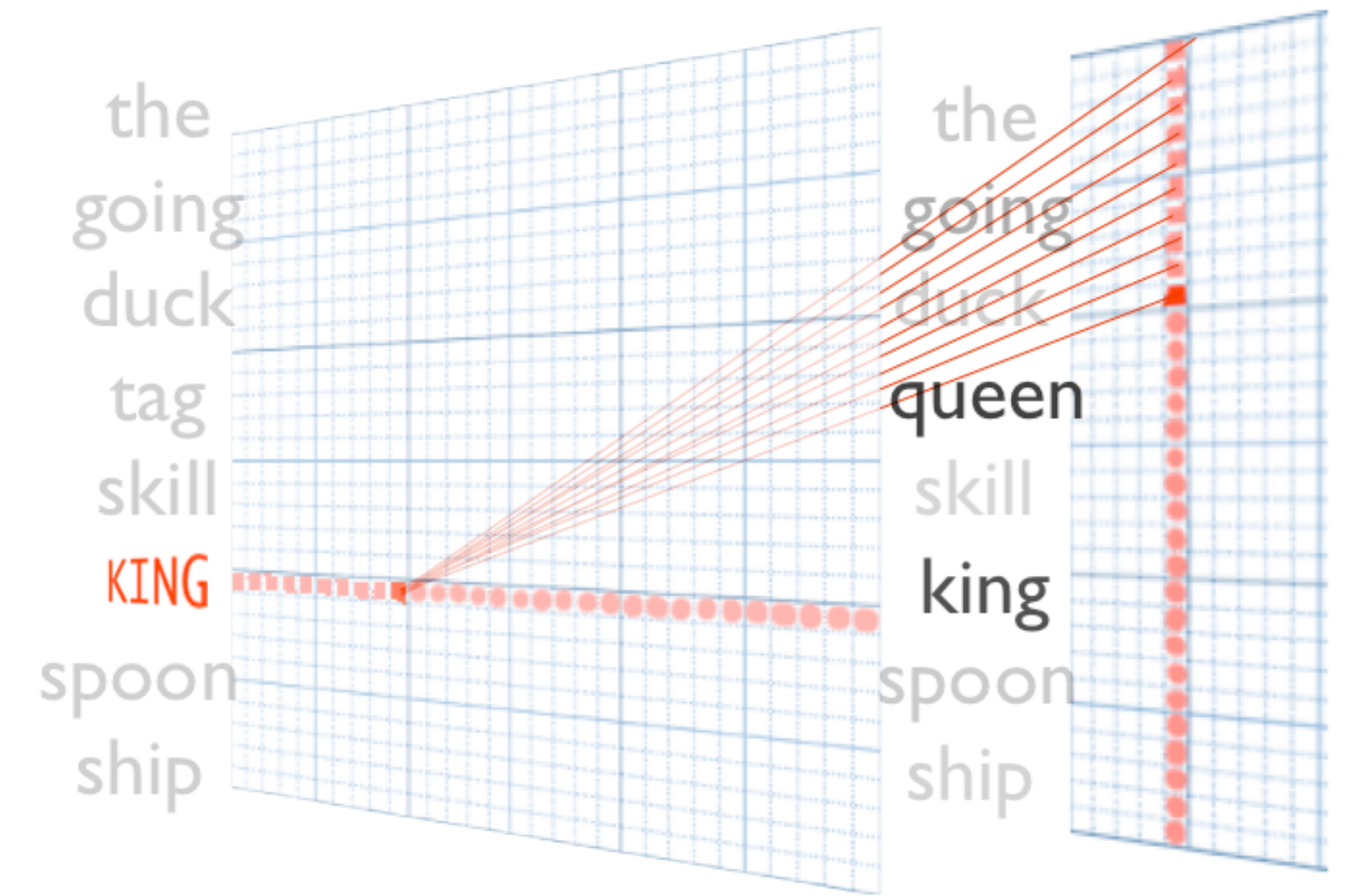
Bag of Words - Curse of Dimensionality

- Before word embeddings - Bag of words
 - Dictionary of words & counts in the text
 - Easy feature generation technique
- Limitations
 - Hard to capture order of words
 - Curse of dimensionality - limited vocabulary - similar words don't match



Word Embeddings Cont'd

- Word Embedding: mapping words to a higher dimensional space, typically 200-500, e.g.,
 - $W(\text{'King'}) = (0.2, -0.4, 0.9, \dots)$
 - $W(\text{'Queen'}) = (0.1, -0.3, 0.8, \dots)$
- Learn representations of words
- How: two layer NN to learn word representations by predicting validity of phrases

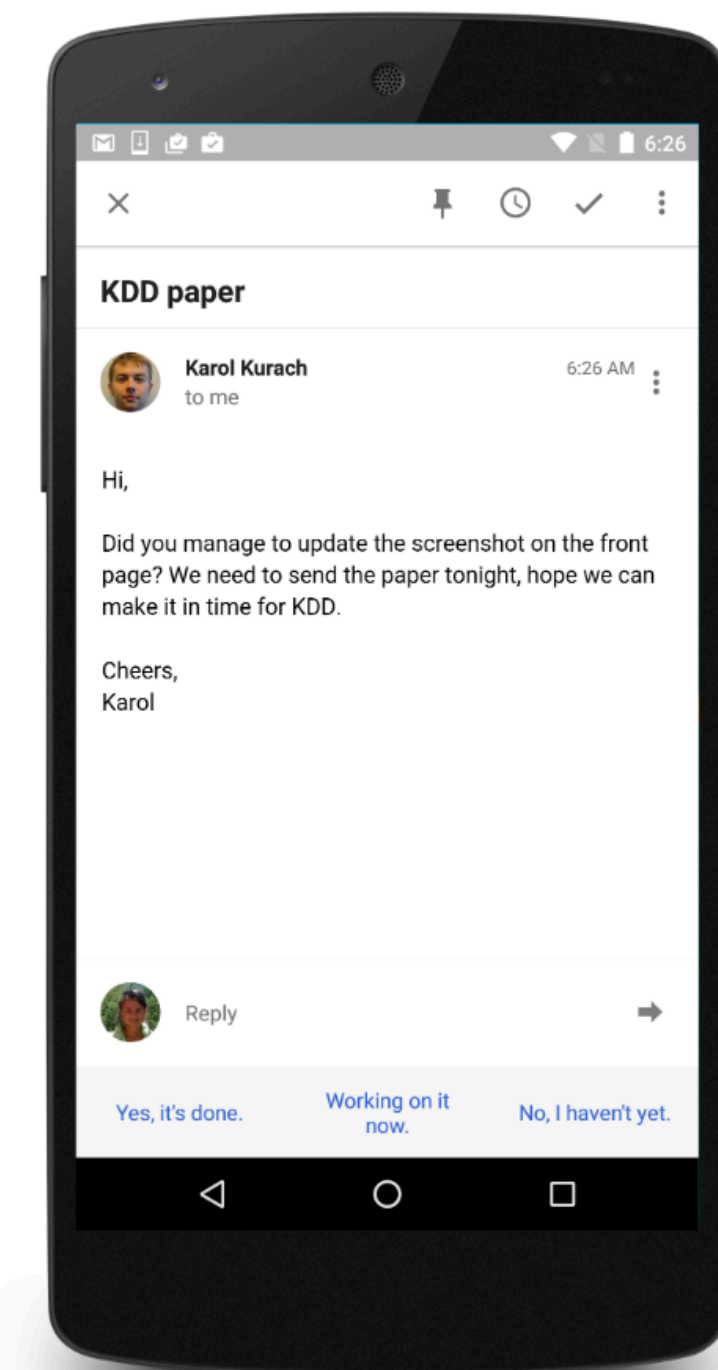
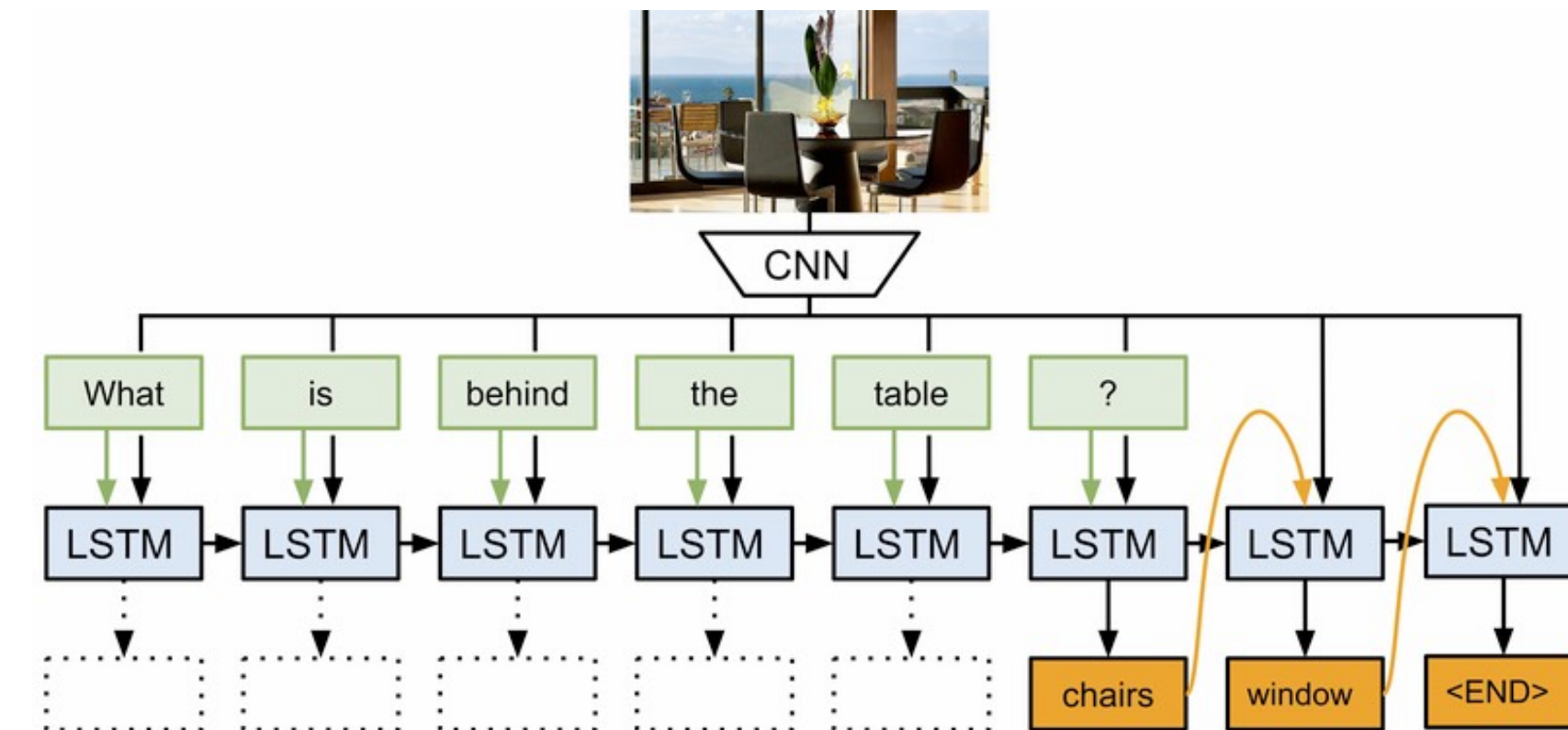


Relationship	Example 1	Example 2	Example 3
France - Paris	Italy: Rome	Japan: Tokyo	Florida: Tallahassee
big - bigger	small: larger	cold: colder	quick: quicker
Miami - Florida	Baltimore: Maryland	Dallas: Texas	Kona: Hawaii
Einstein - scientist	Messi: midfielder	Mozart: violinist	Picasso: painter
Sarkozy - France	Berlusconi: Italy	Merkel: Germany	Koizumi: Japan
copper - Cu	zinc: Zn	gold: Au	uranium: plutonium
Berlusconi - Silvio	Sarkozy: Nicolas	Putin: Medvedev	Obama: Barack
Microsoft - Windows	Google: Android	IBM: Linux	Apple: iPhone
Microsoft - Ballmer	Google: Yahoo	IBM: McNealy	Apple: Jobs
Japan - sushi	Germany: bratwurst	France: tapas	USA: pizza

Example of similar word vectors

Sequence Learning: Response Generation

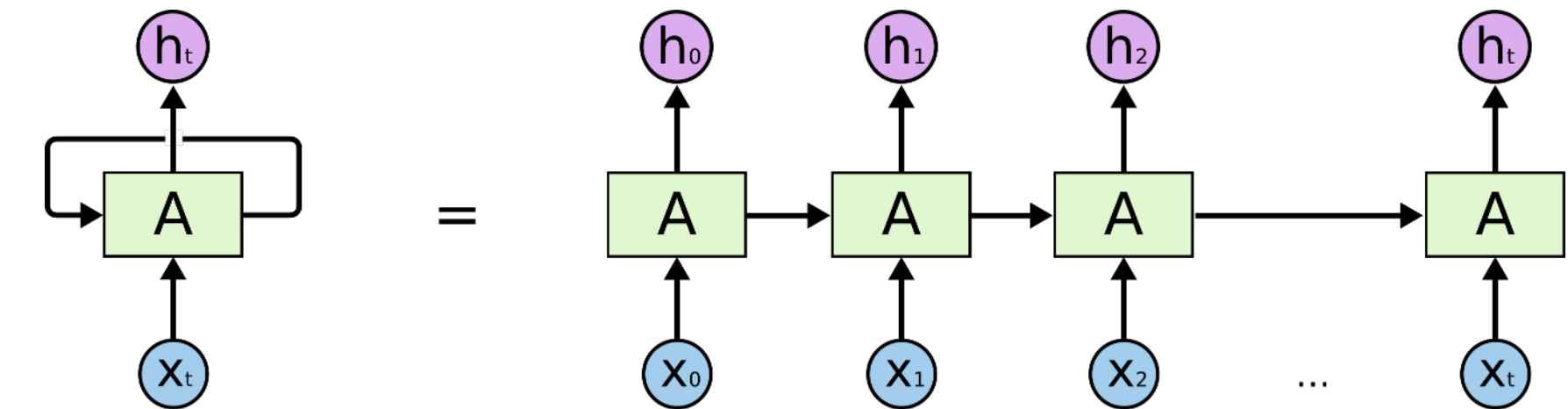
- Automated Response Generation
- Sequence 2 Sequence Model
- Recurrent Neural Network (RNN)
- Long Short Term Memory Network (LSTM)
- Example: GMail Smart Reply
- Automated Response Suggestions



Sequence Learning: RNN And LSTM

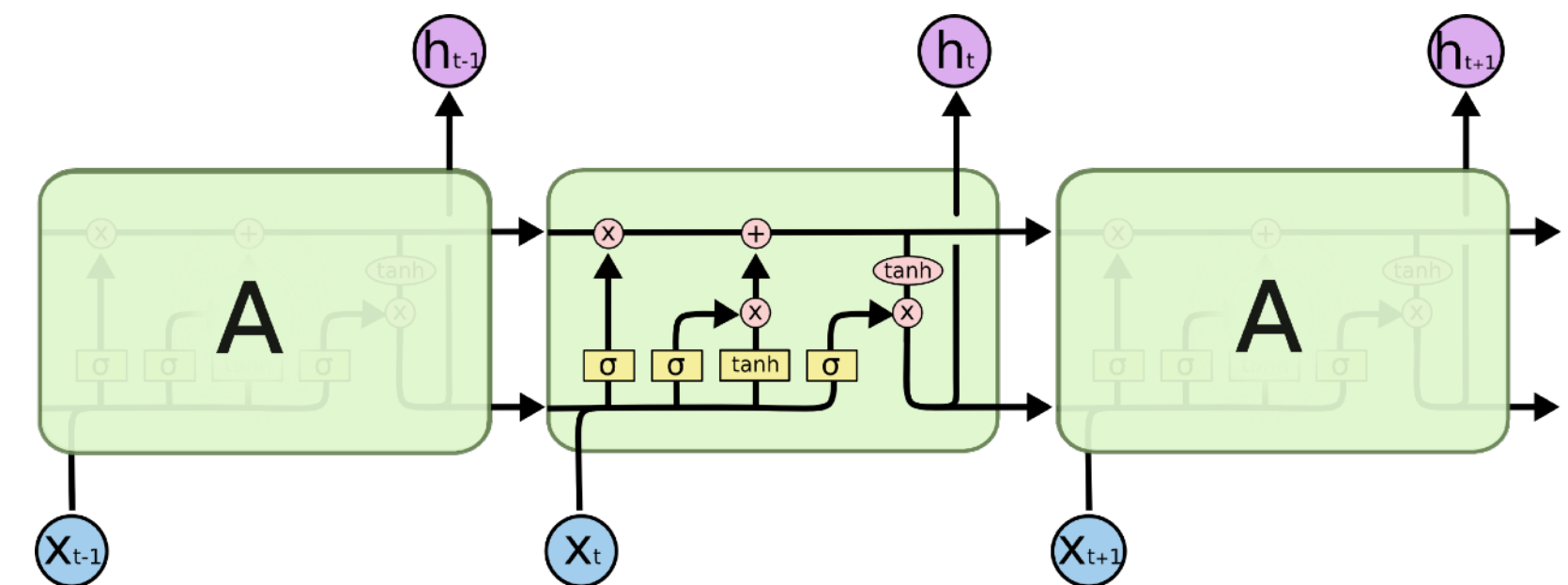
- Recurrent Neural Network

- Output of a module go into a module of same type (recurrent)
- Good for capturing a sequence



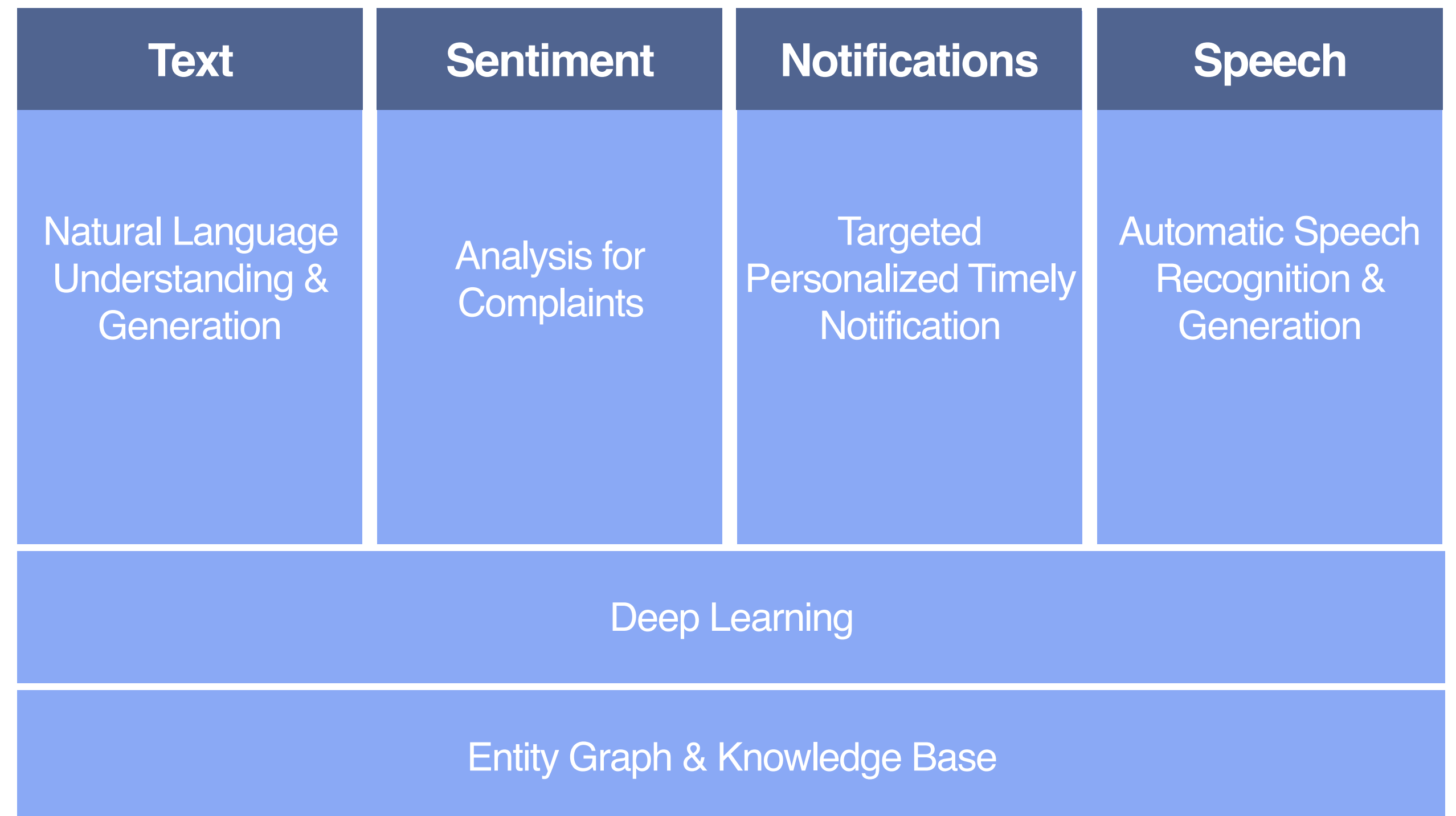
- Long Short Term Memory Network

- Long running cell state: forget & add new values
- Output: combination of cell state, previous output, and new input



Training Deep Learning Models for NLP

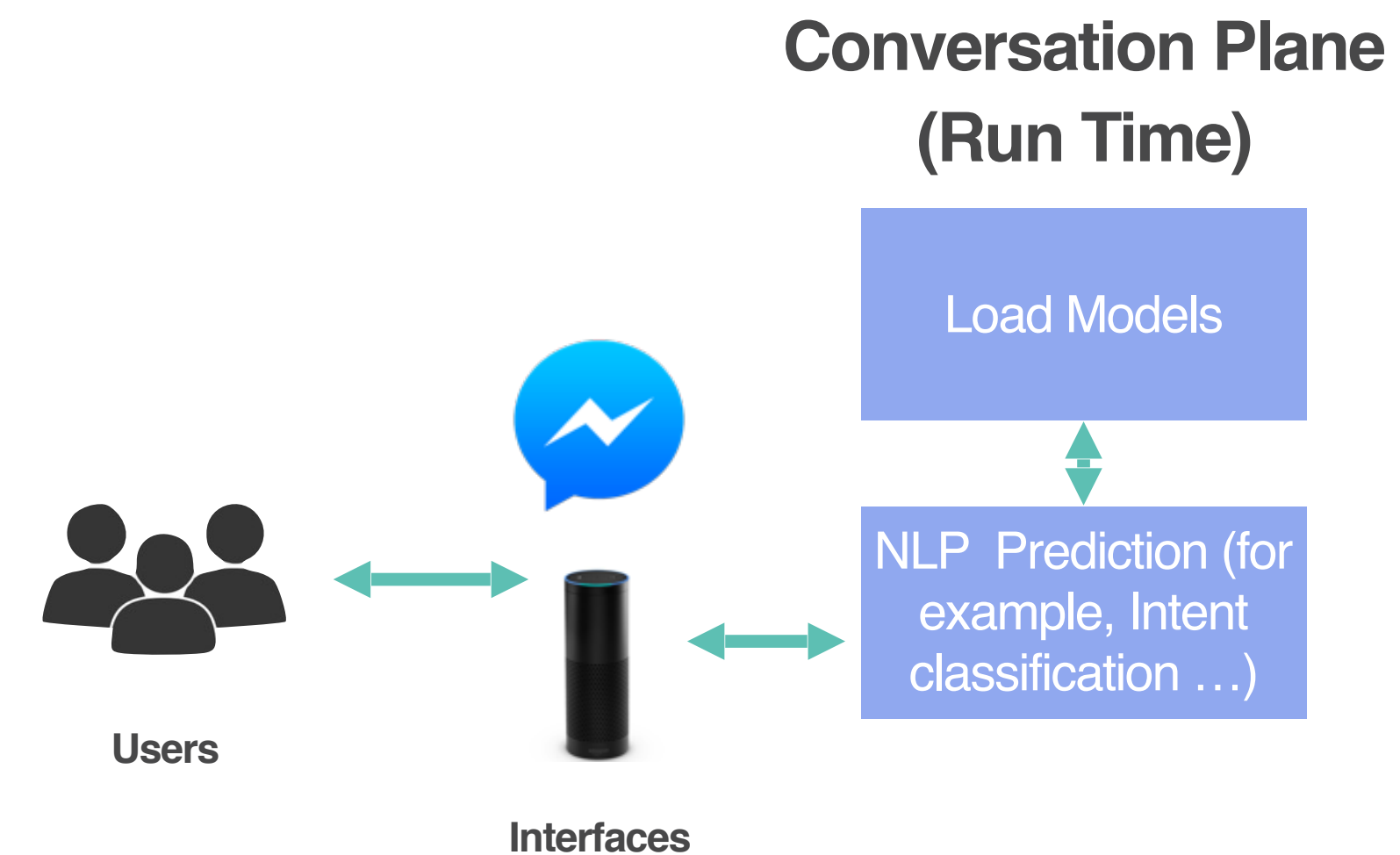
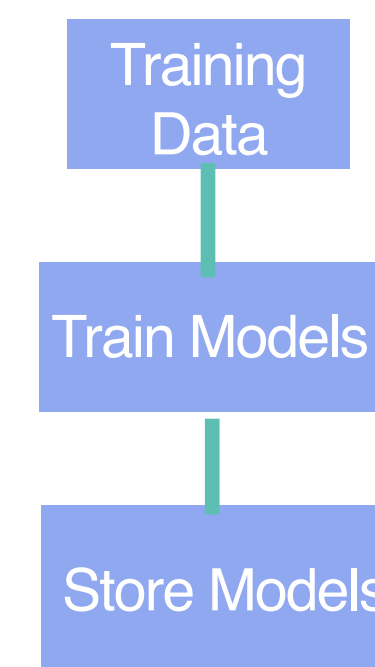
- Intent Classification
 - Deep Learning — LSTM
- Information Extraction
 - Named Entity Recognition (NER)
 - Slot attributes
- Sentiment and Complaint Classifier
- Knowledge Base & Semantic Search
- Machine Reading Comprehension



Scaling Training Deep Learning Models For NLP

- Off line: Started with a script for training models
- Run Time: A service for prediction during runtime
- However, the number of models are reaching in thousands
- Hard to manage model training script for each of the bot

Control Plane (Offline)



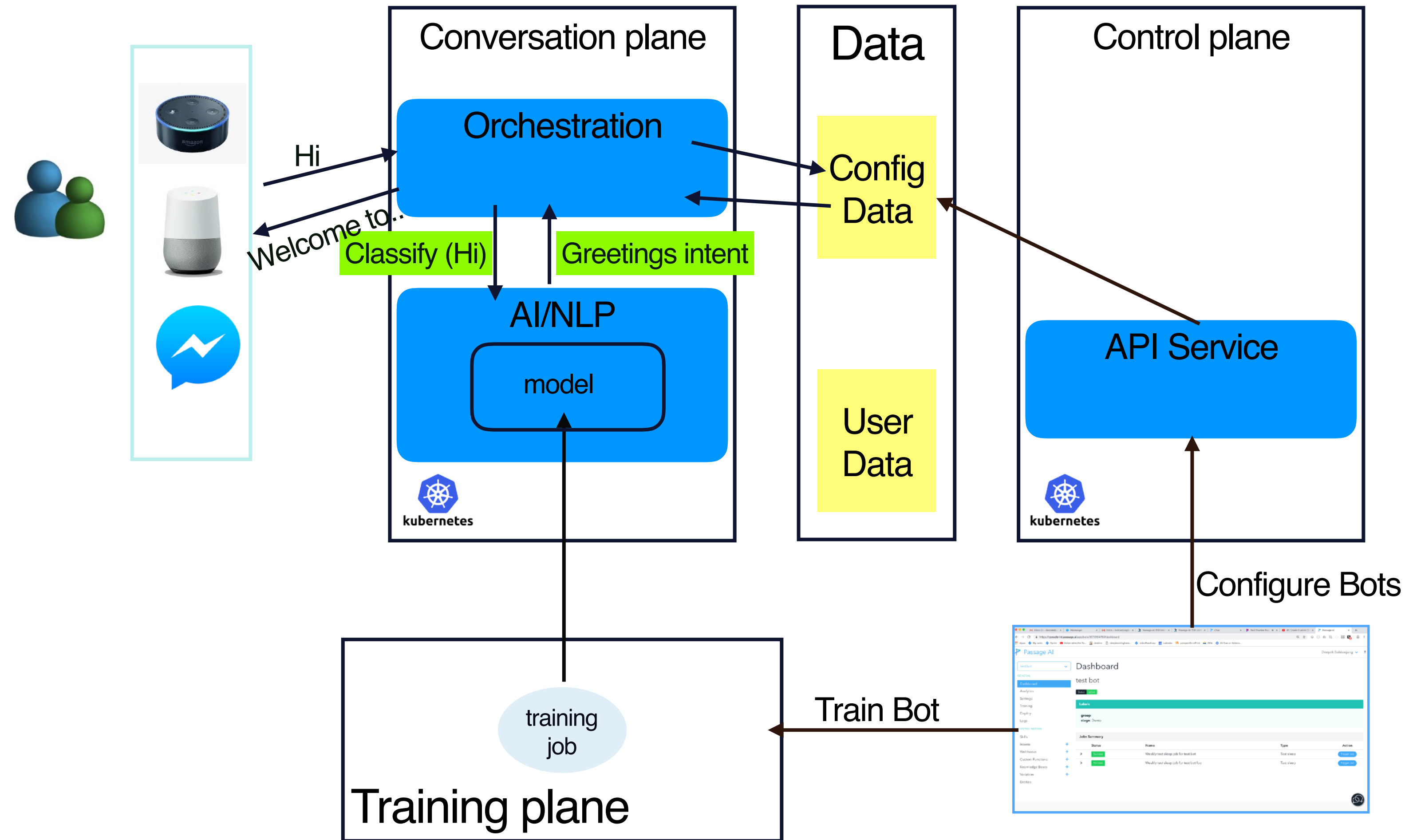
Outline

Conversational AI and Deep Learning

Need for a Jobs framework on Kubernetes

Our Jobs architecture

Passage AI Architecture



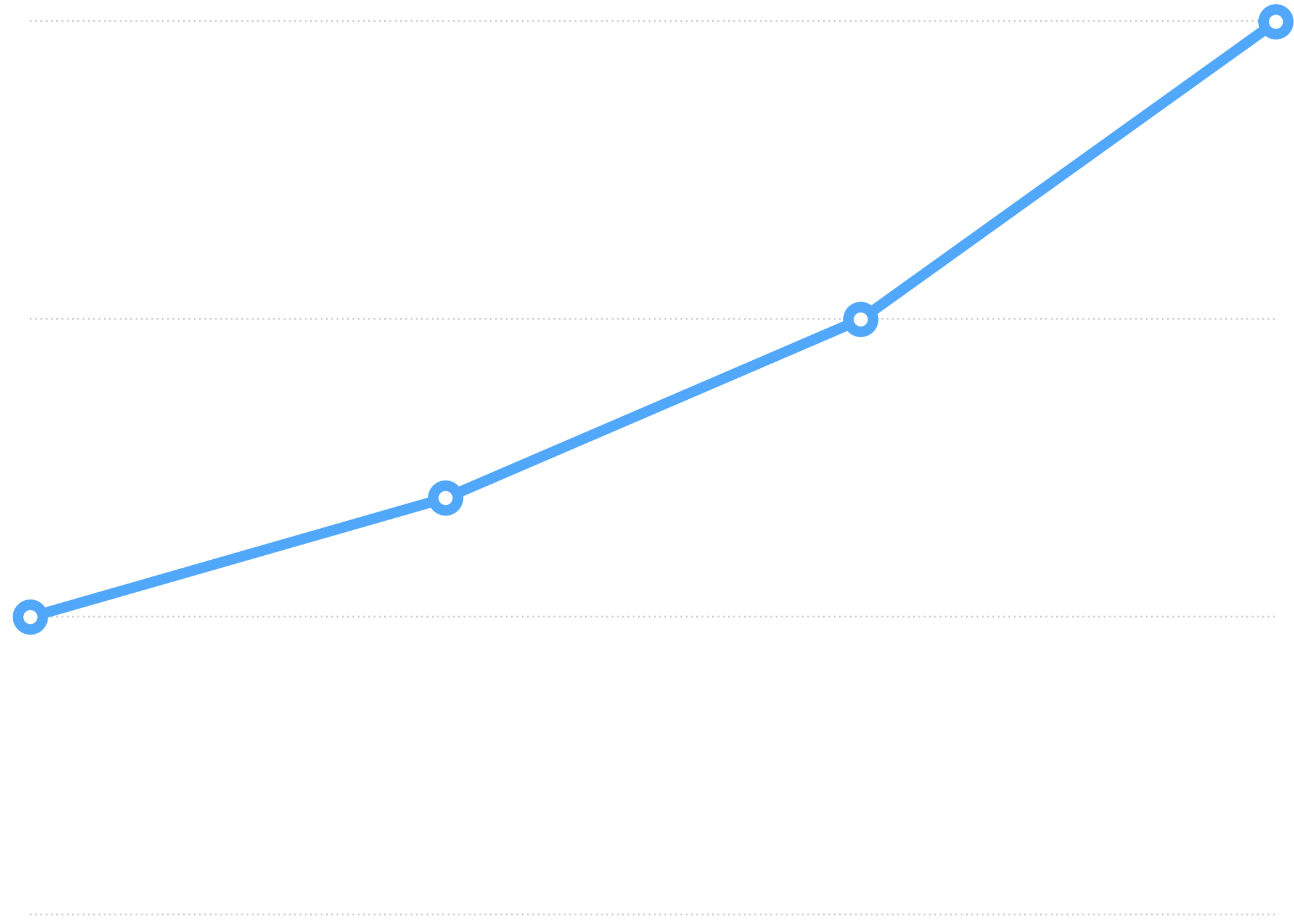
Passage AI Architecture



When do we train a new model for a bot ?

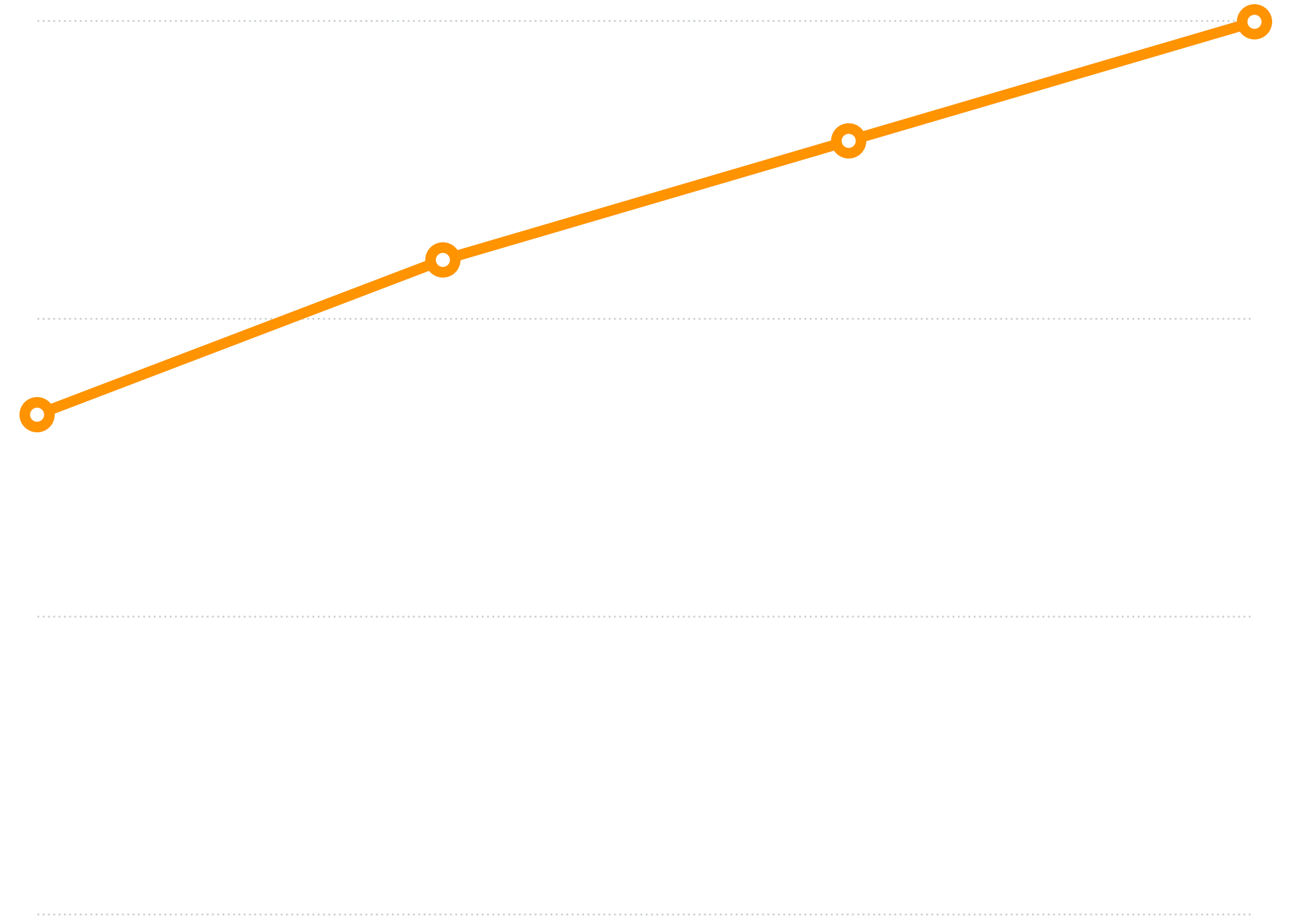
- When a new bot is created
- When a bot is changed
 - utterances are added or modified
 - New training data is available

of Bots



August September October November

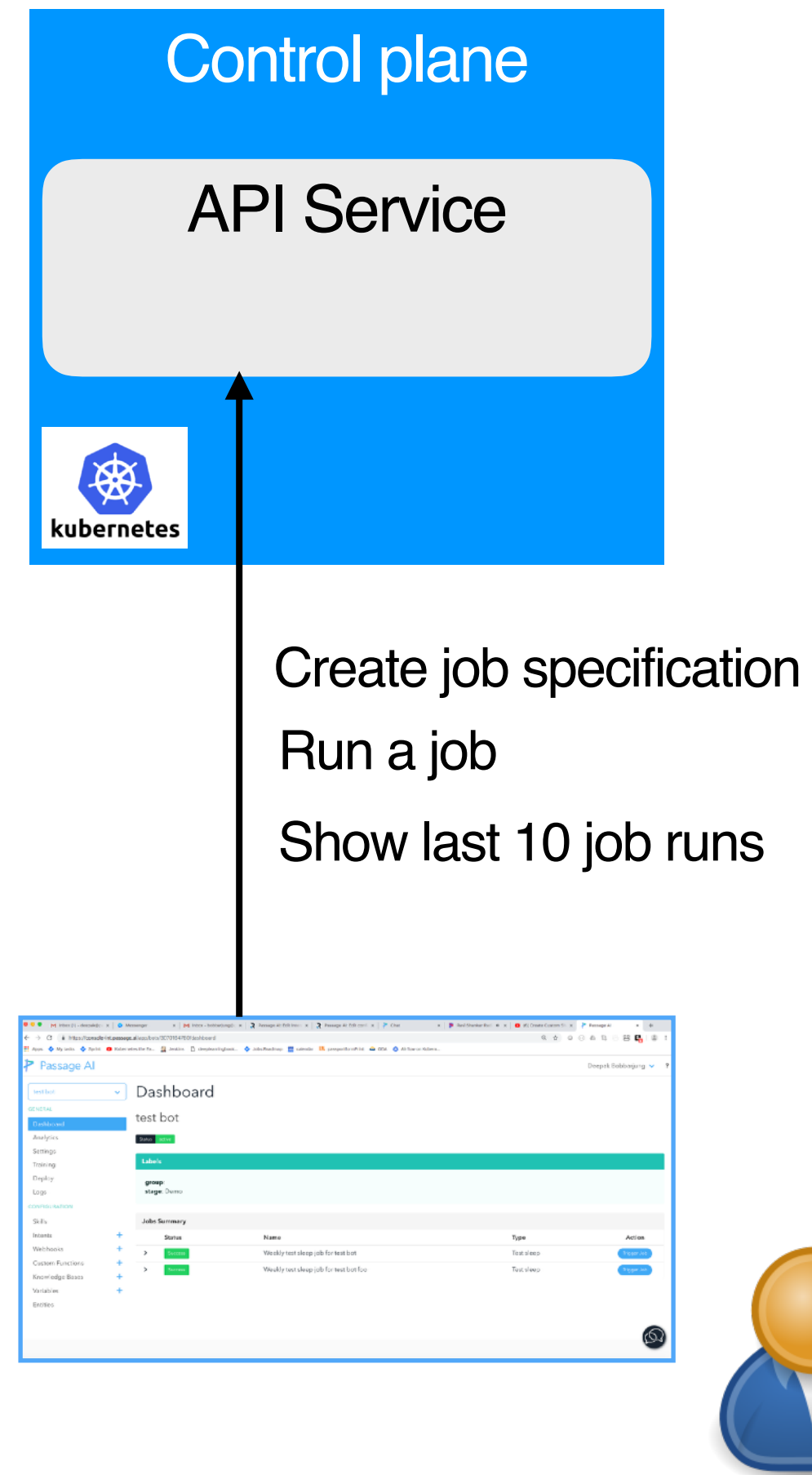
of Bot changes per day



August September October November

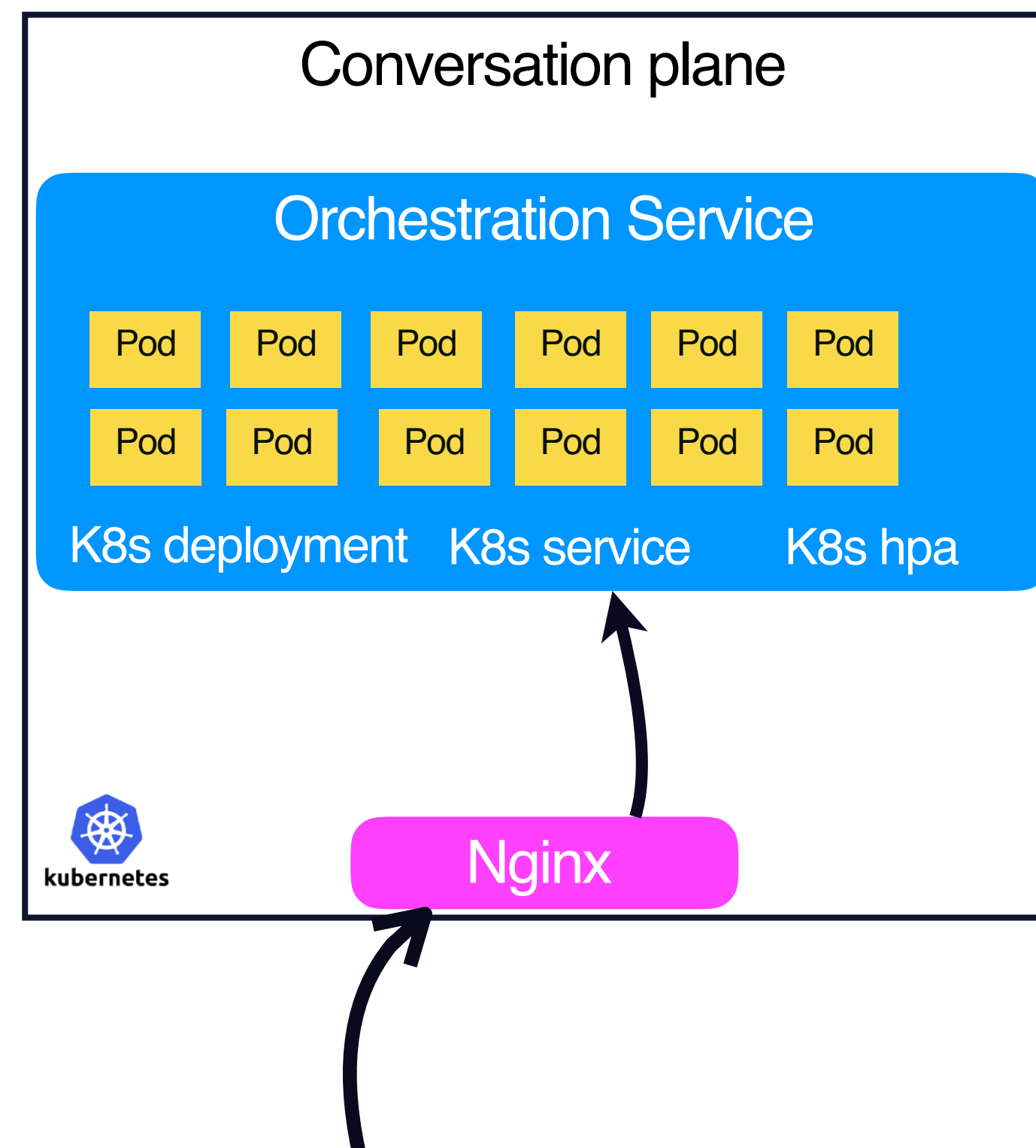
Why do we need a Jobs framework?

- Run jobs at scale
- Eliminate out of band scripts that tend to become 'tribal'.
- APIs and UI for exposing jobs to our customers in our Bot Builder UI.
- Reporting and auditing around jobs.



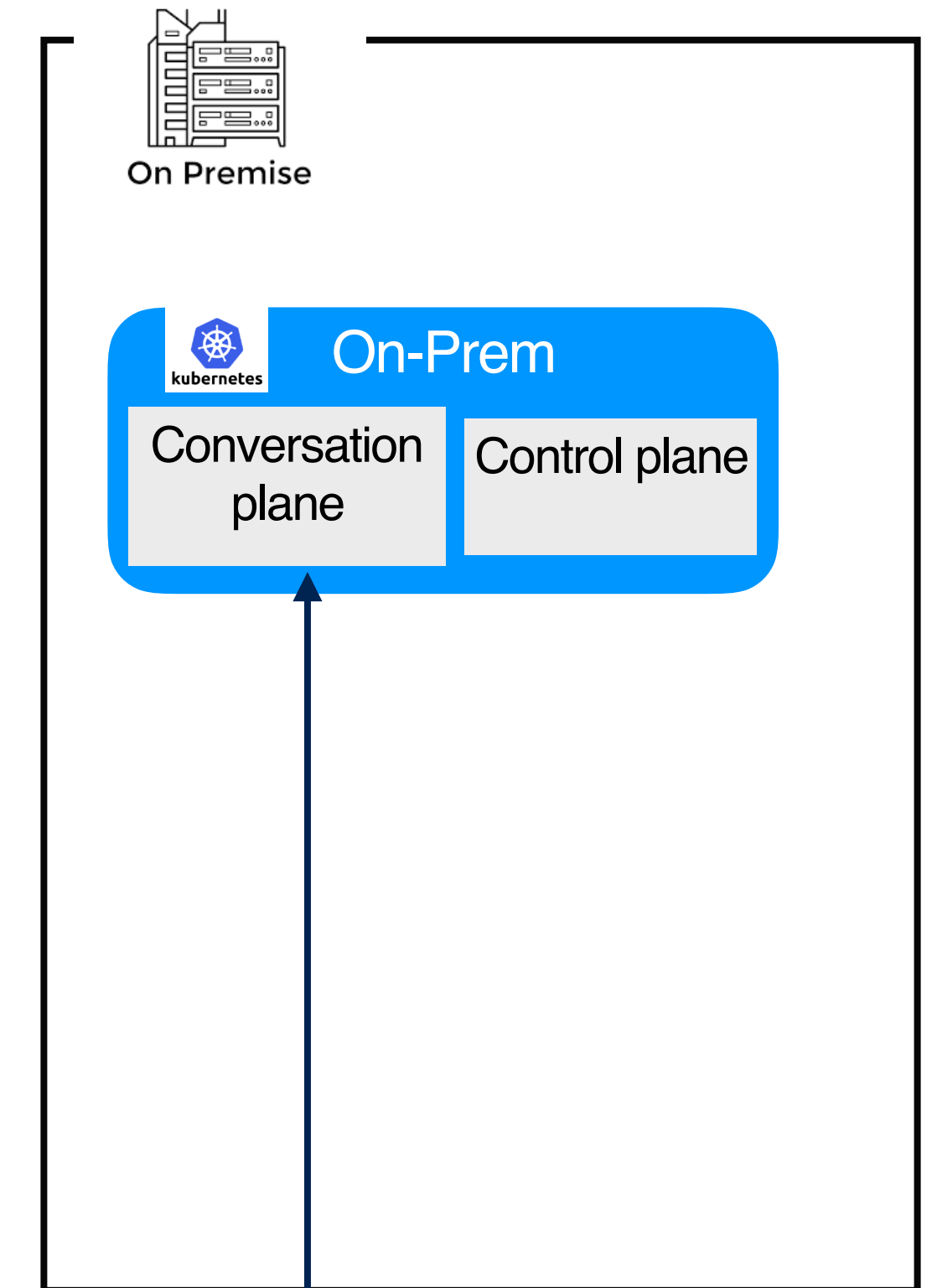
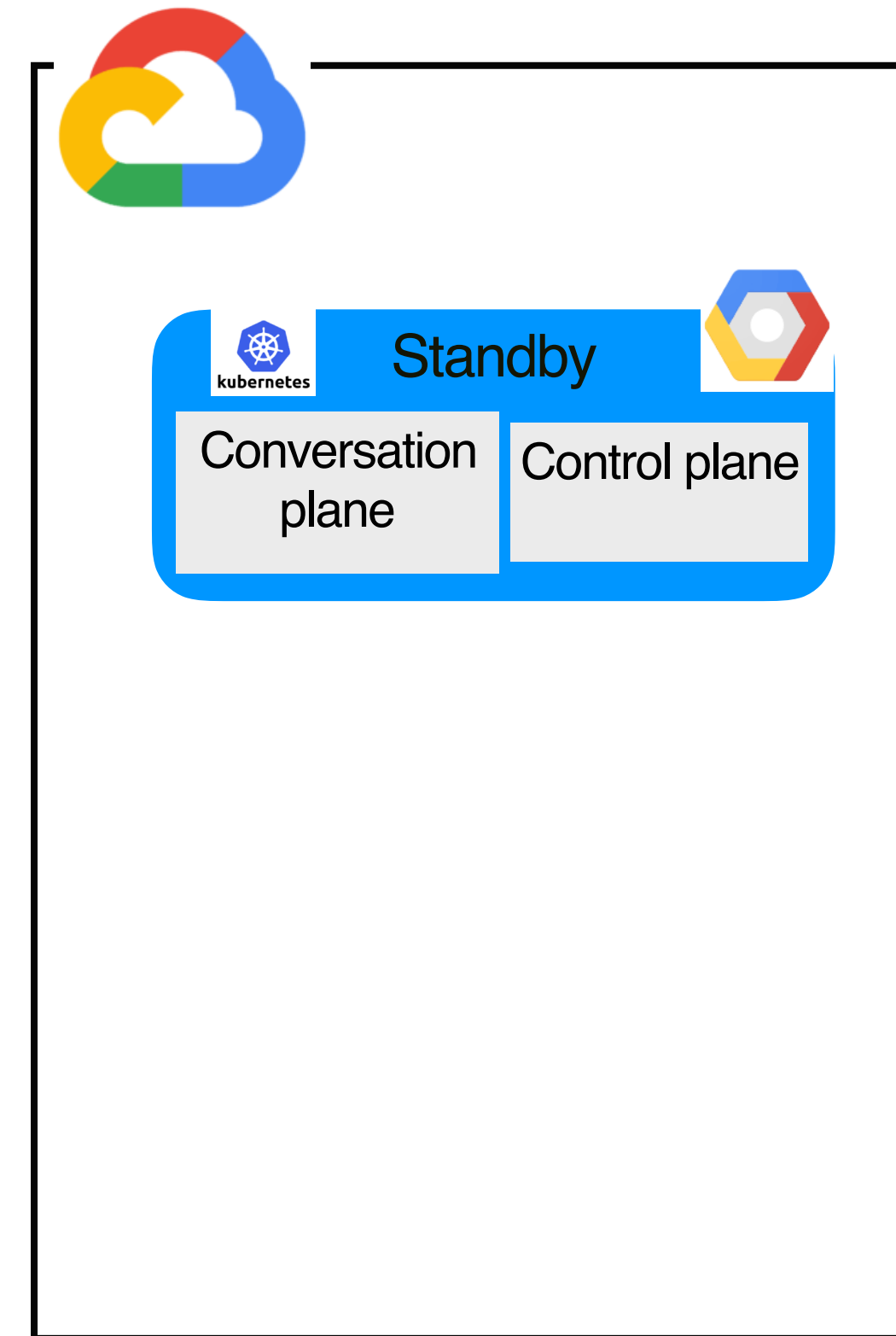
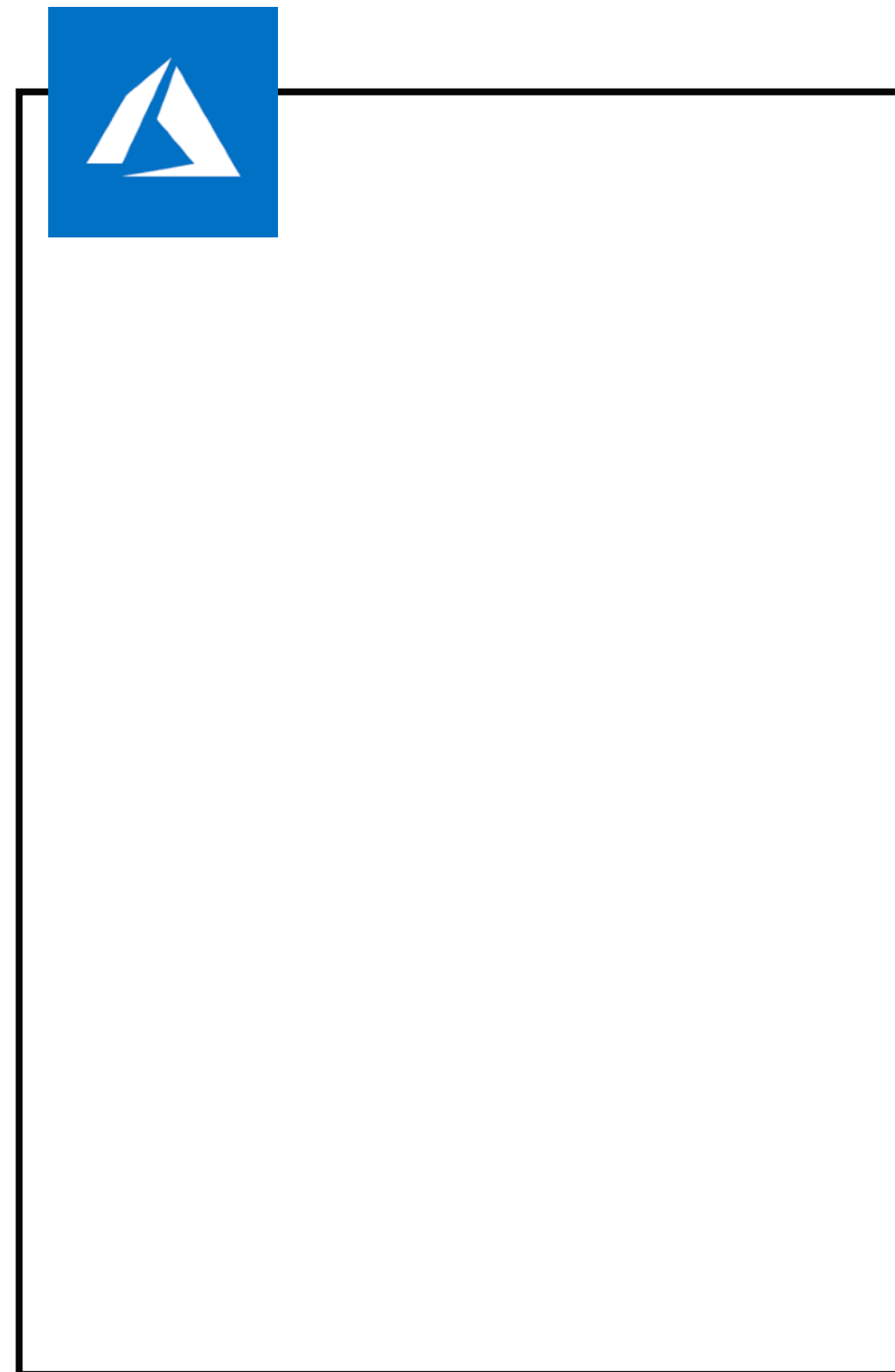
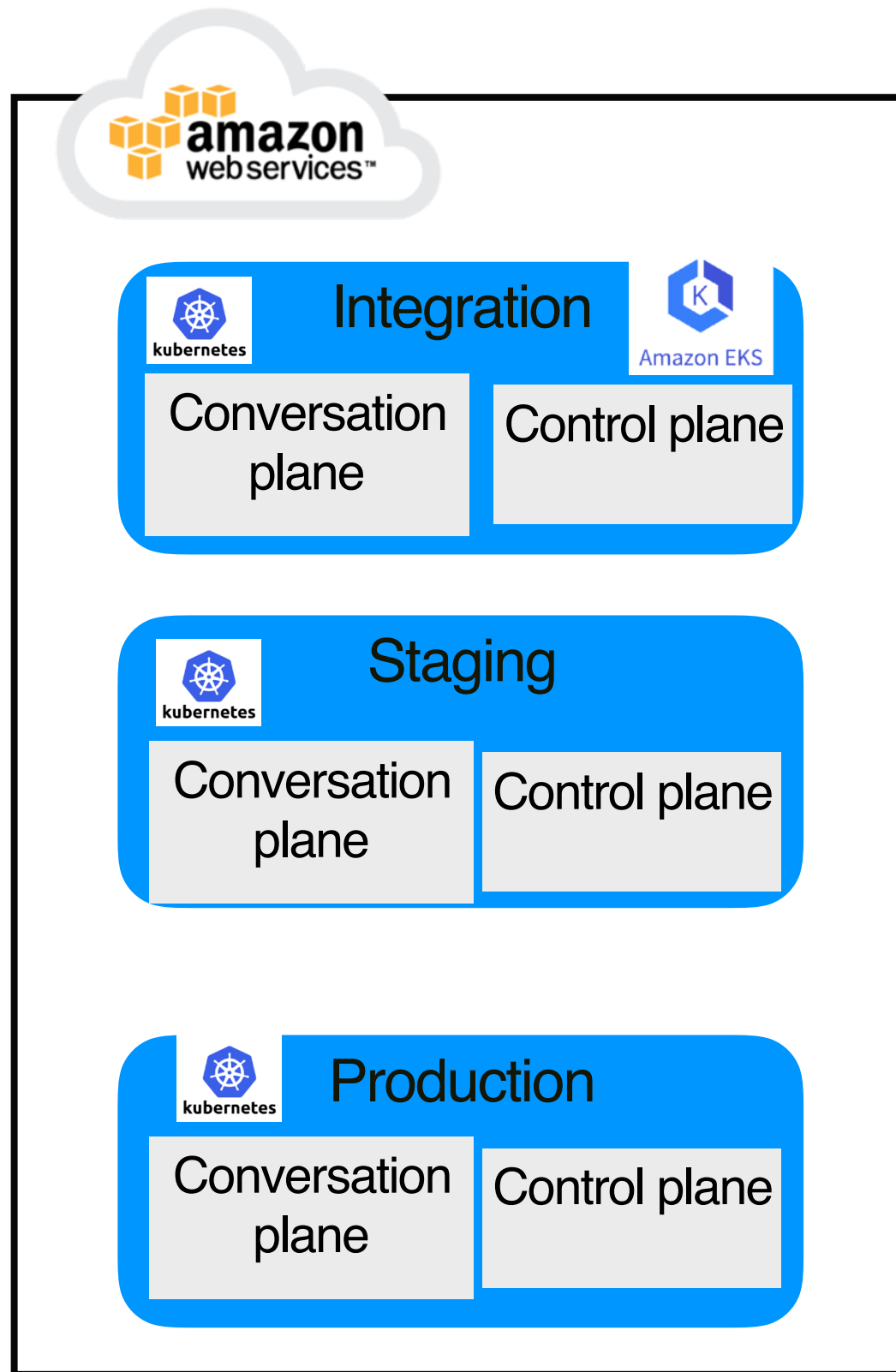
Why Kubernetes (K8S) For Our Microservices?

- Scale and availability of our microservices



Why Kubernetes?(Contd)

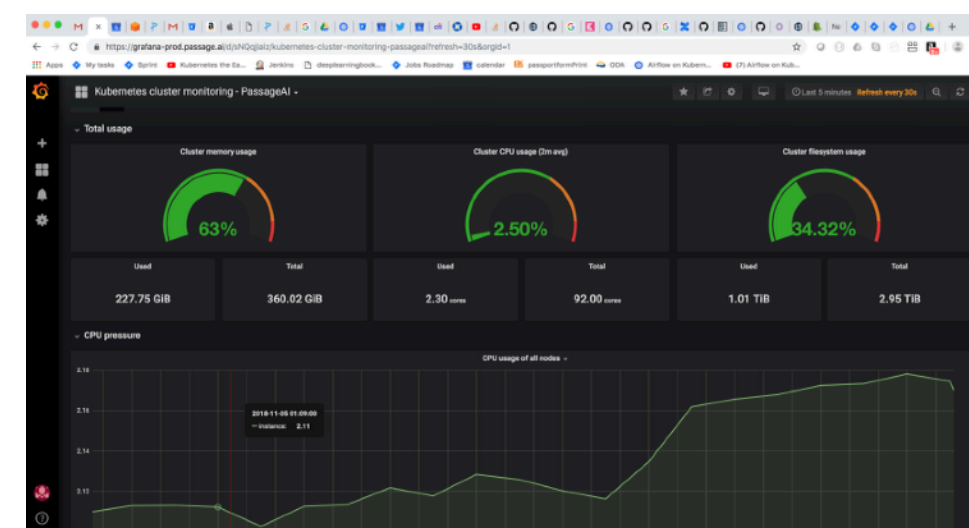
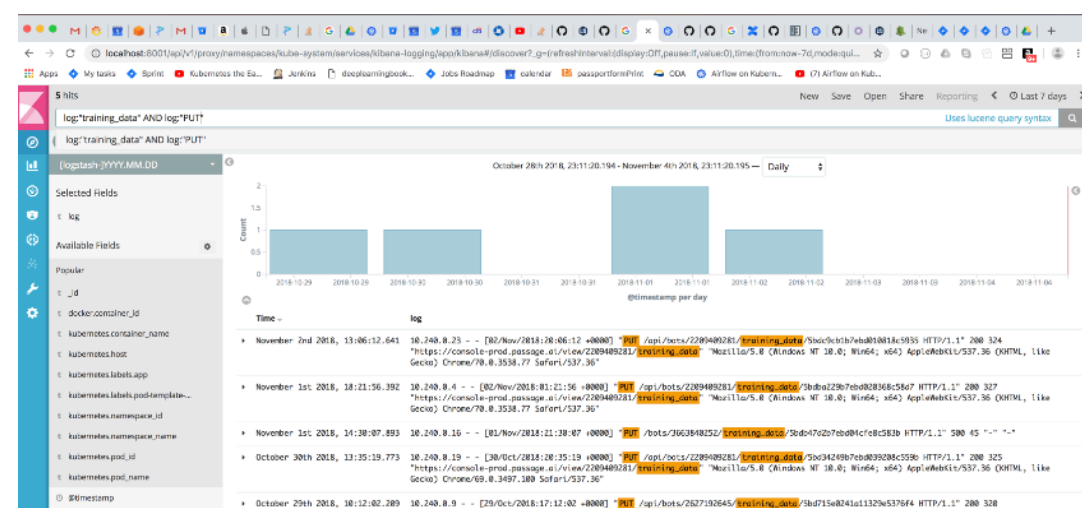
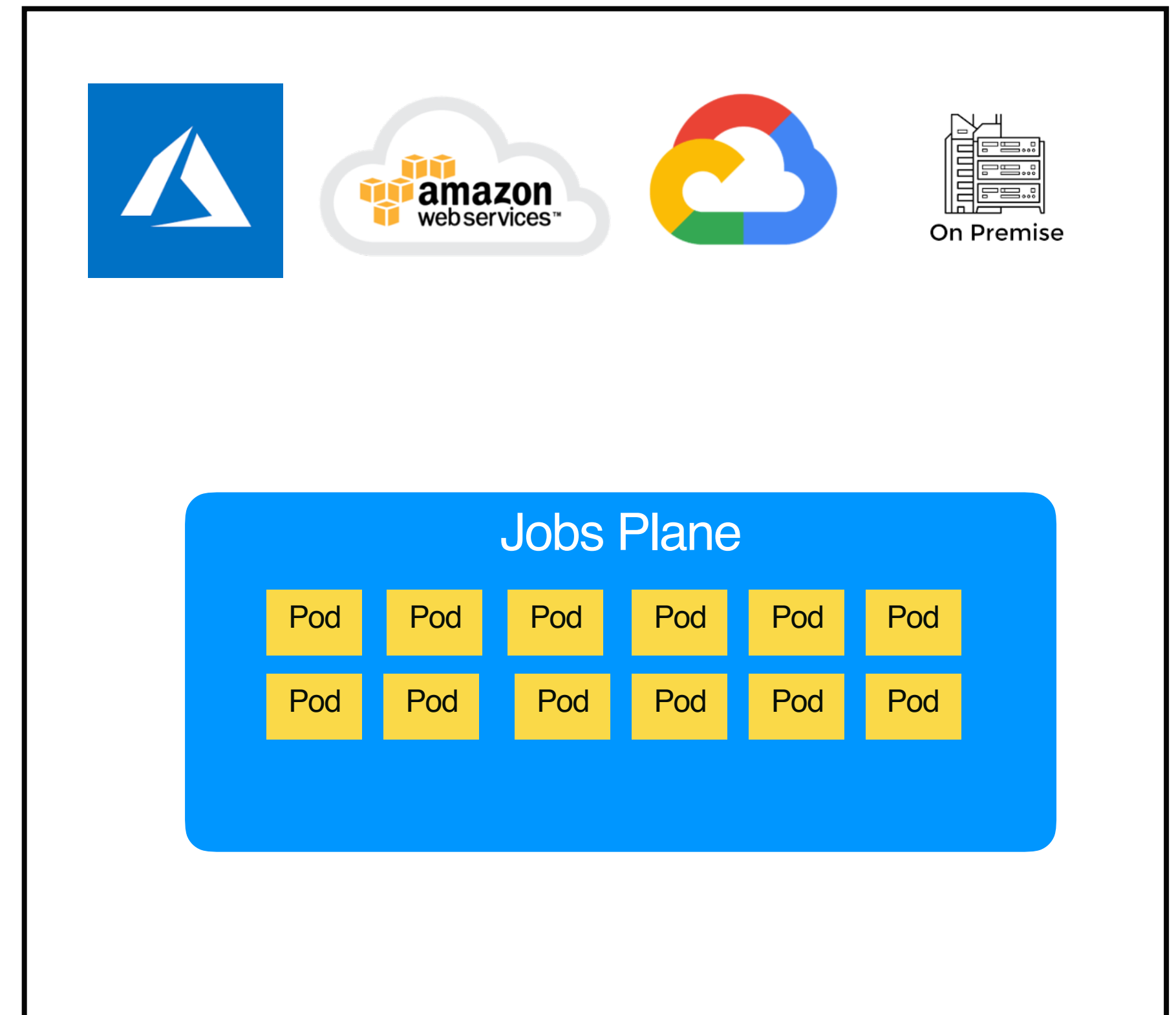
- Cloud Agnostic and On-prem ready



\$helm install passage-ai

Why Create The Jobs Framework In Kubernetes?

- Jobs should also be cloud-agnostic and on-prem ready
- Handle scale and availability in the same way as our microservices
- Same set of tools for monitoring, logging and auditing.



Outline

Conversational AI and Deep Learning

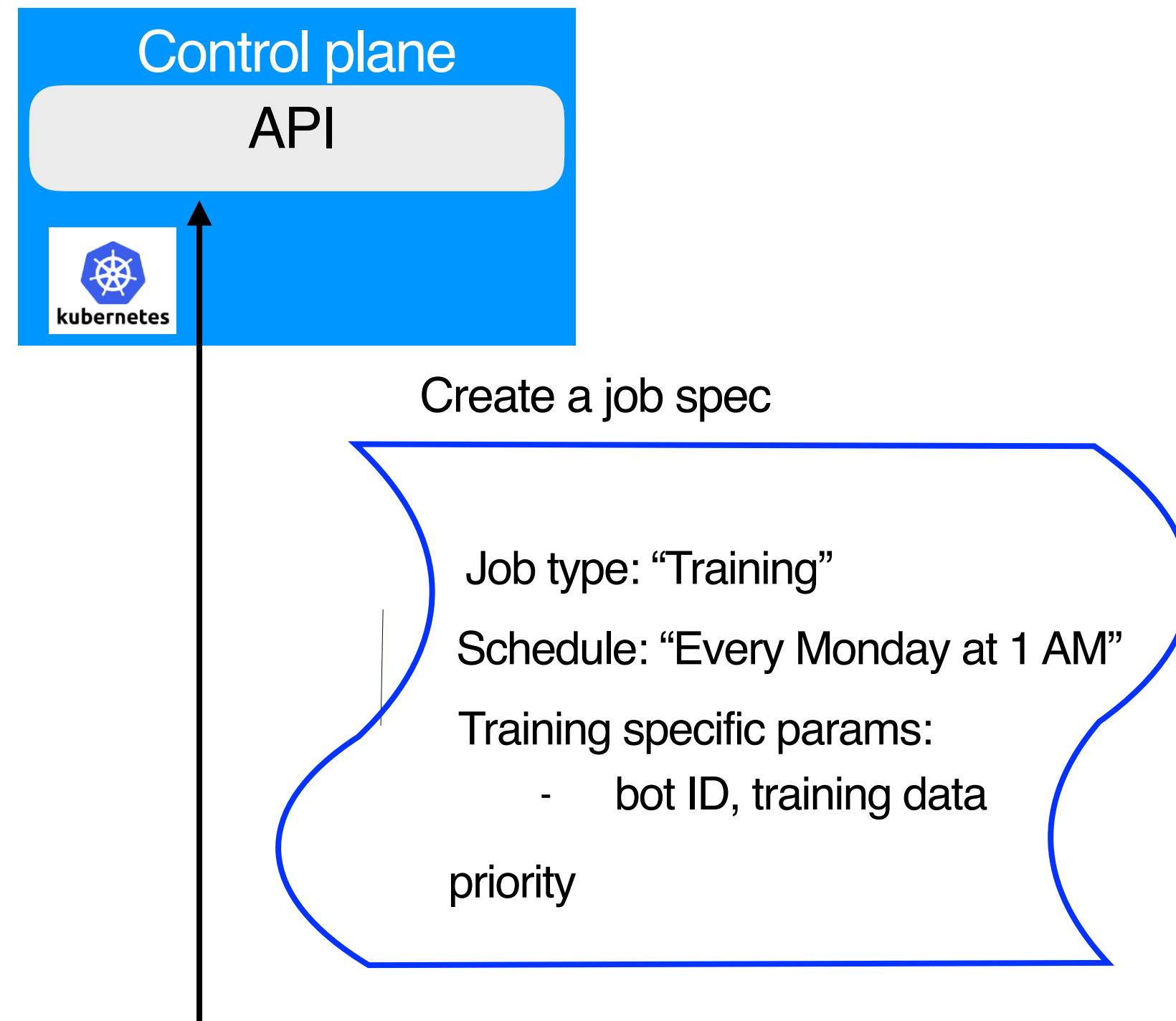
Need for a Jobs framework on Kubernetes

Our Jobs architecture

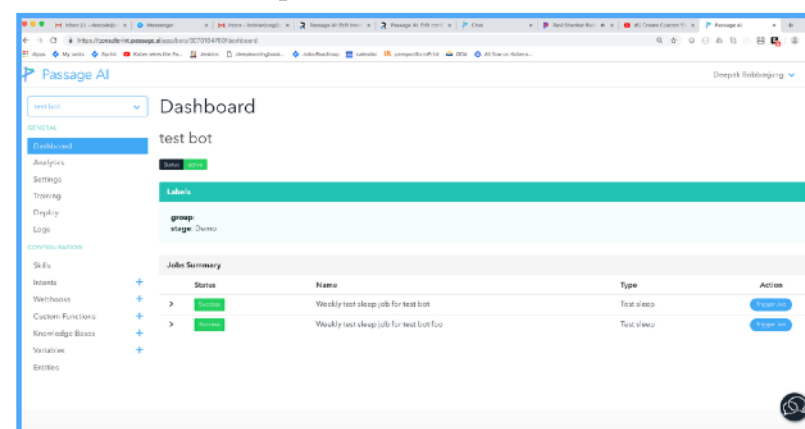
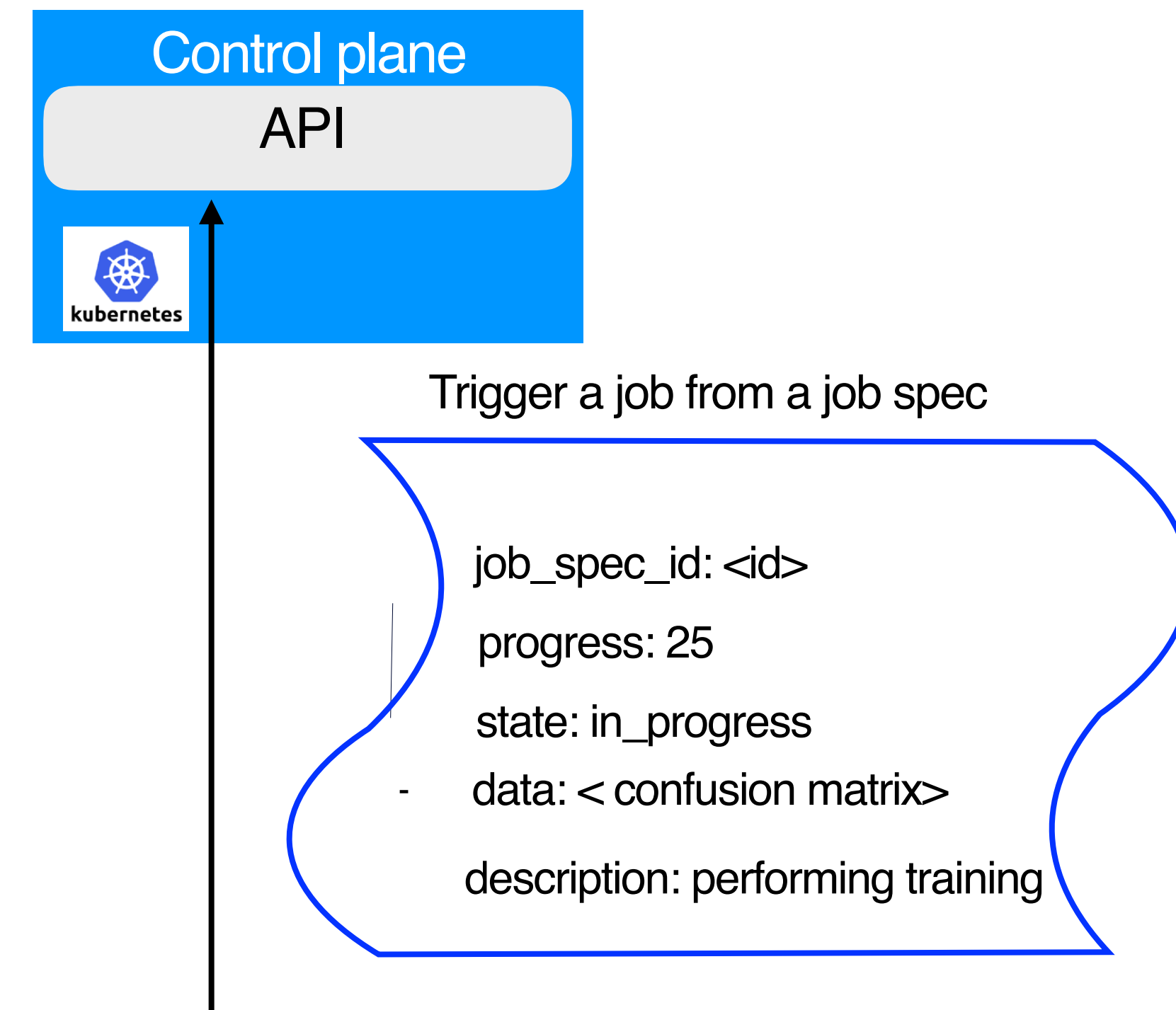
Example Job types in our system

- Training deep learning models
- Extracting and indexing knowledge base articles
- Nightly testing of our bots

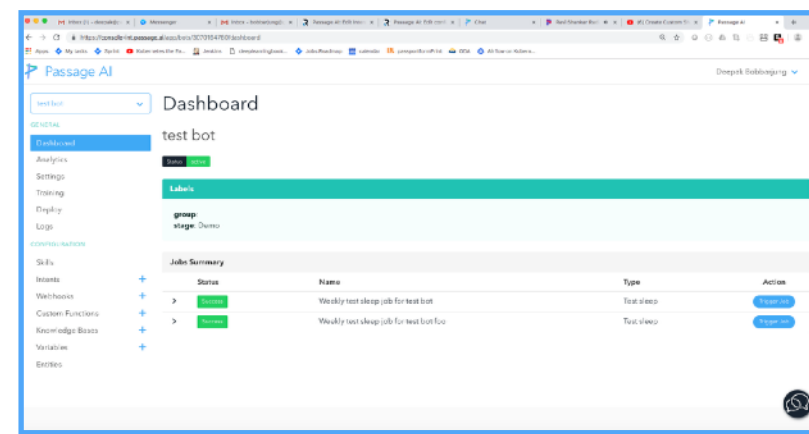
Job Specification



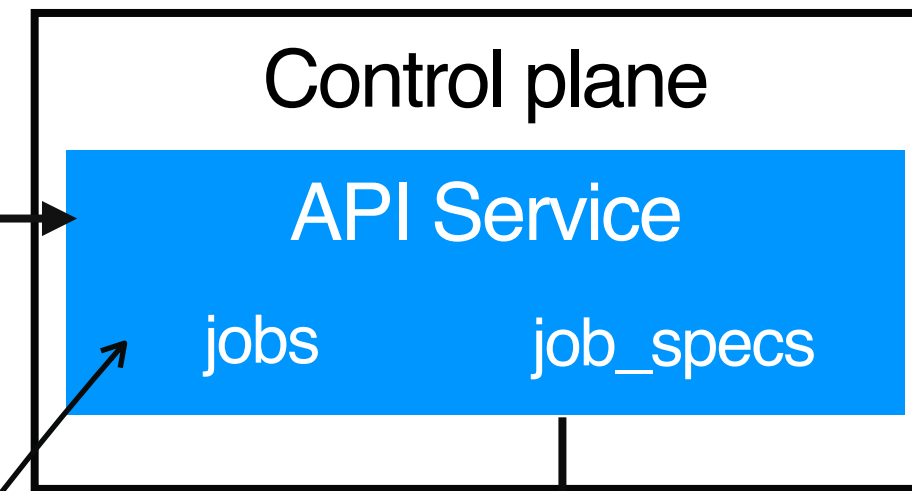
Job Object



Jobs Architecture



Trigger training job from job spec



Jobs plane

Create a Job (params)

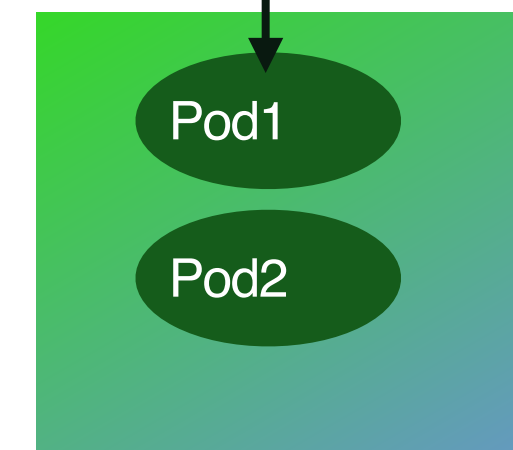
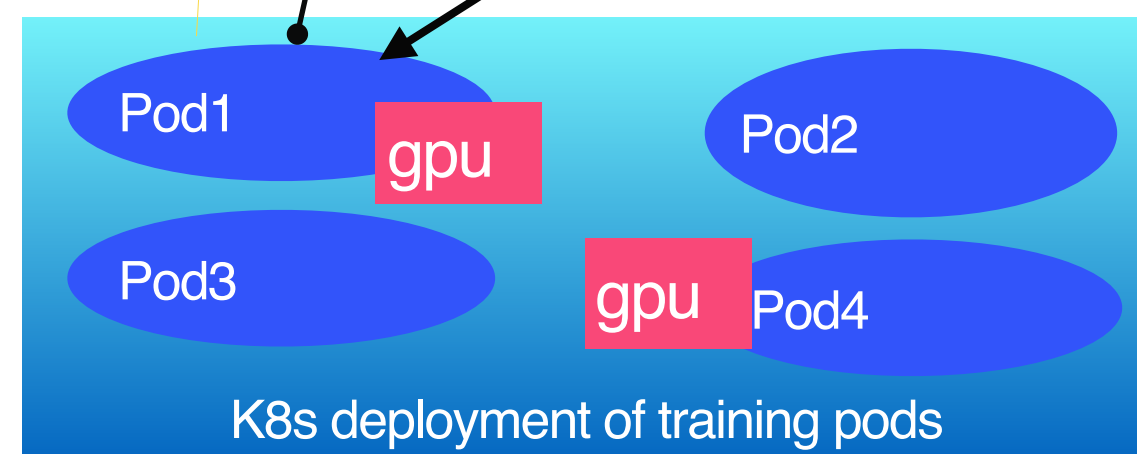


Update job progress

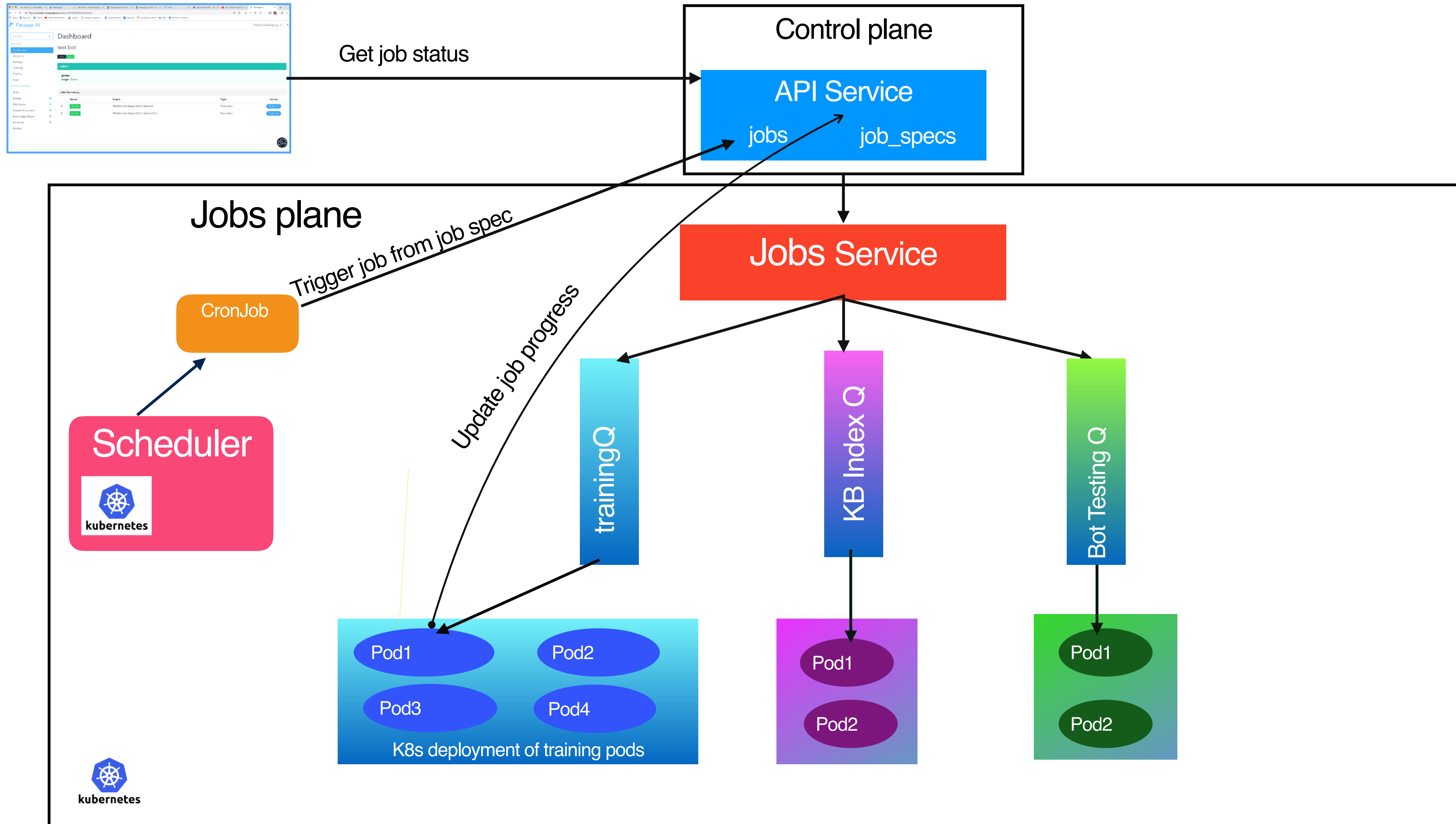
Add item on Q



Pickup item



Jobs Architecture (scheduled jobs)



Alternatives that we considered



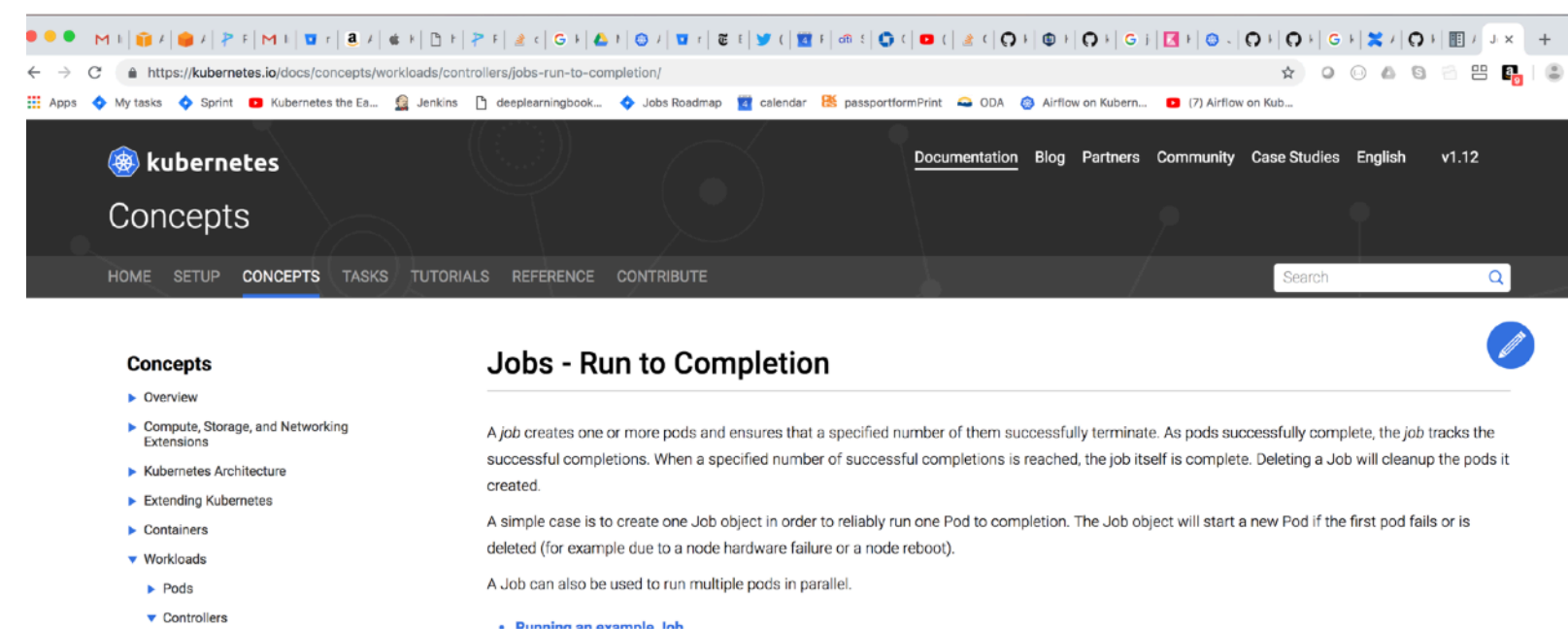
Apache Airflow



Kubeflow



Azkaban



Thank You

 Passage AI

Mitul Tiwari

mitul@passage.ai

[@Mitultiwari](https://twitter.com/Mitultiwari)

Deepak Bobbarjung

deepak@passage.ai

[@Bobbarjung_](https://twitter.com/Bobbarjung_)