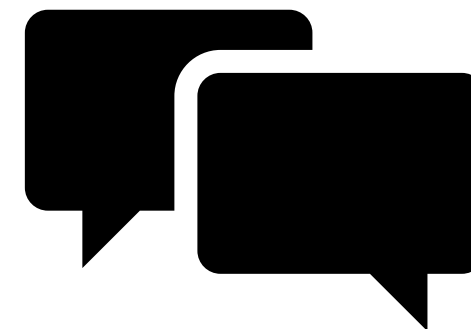
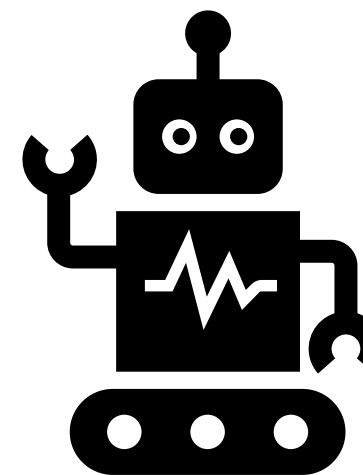


# Demystifying Chatbots



SHIVAM SINHA  
17-11-2019

# What are Chatbots?



Chat – Conversational / series of tasks (message exchange)

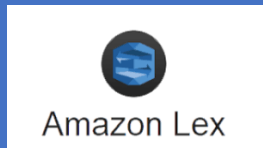


Bots – Cognitive Agents / narrow or Weak AI

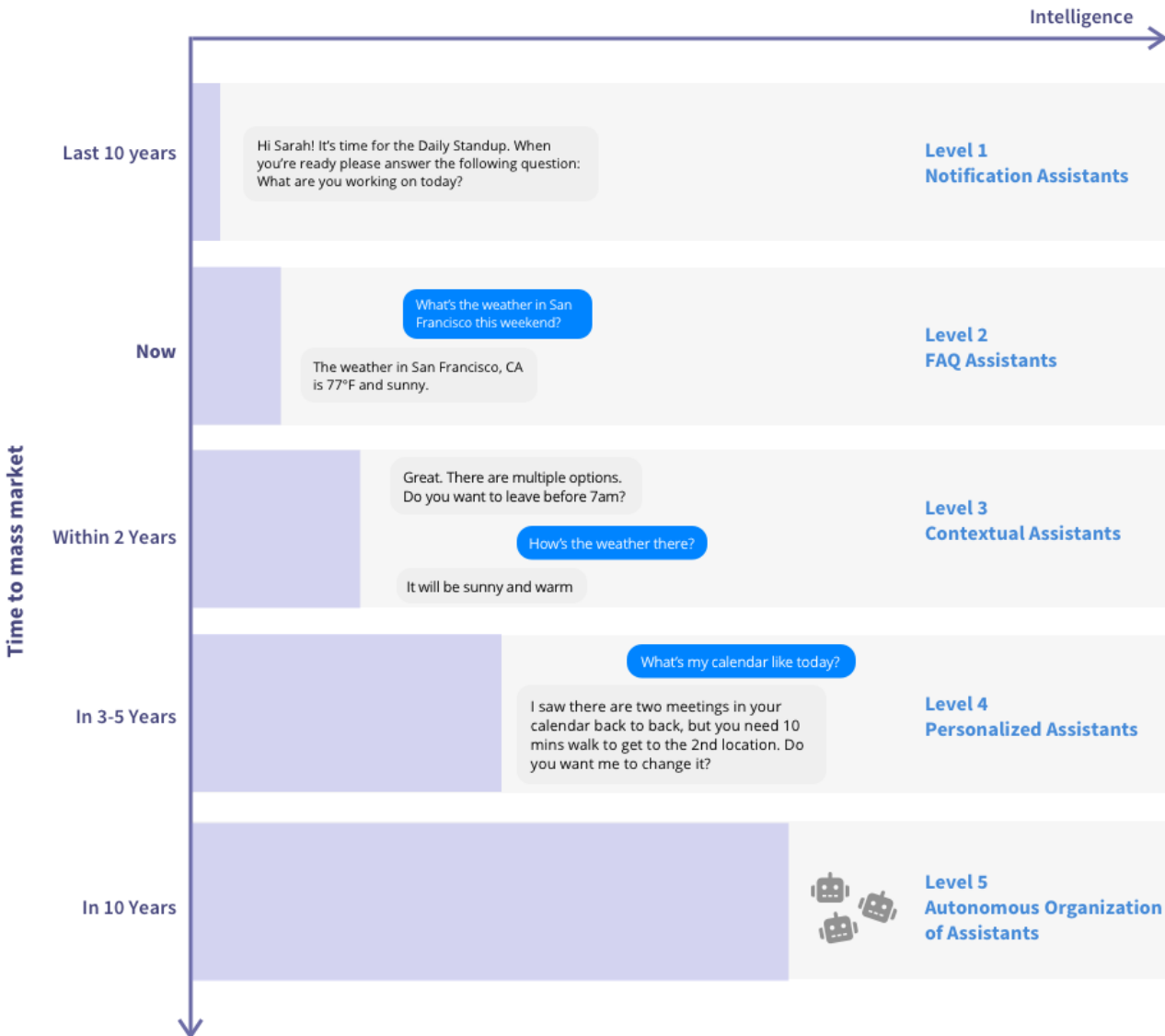
92 million

3.5 billion

80%



# Types of Chatbots



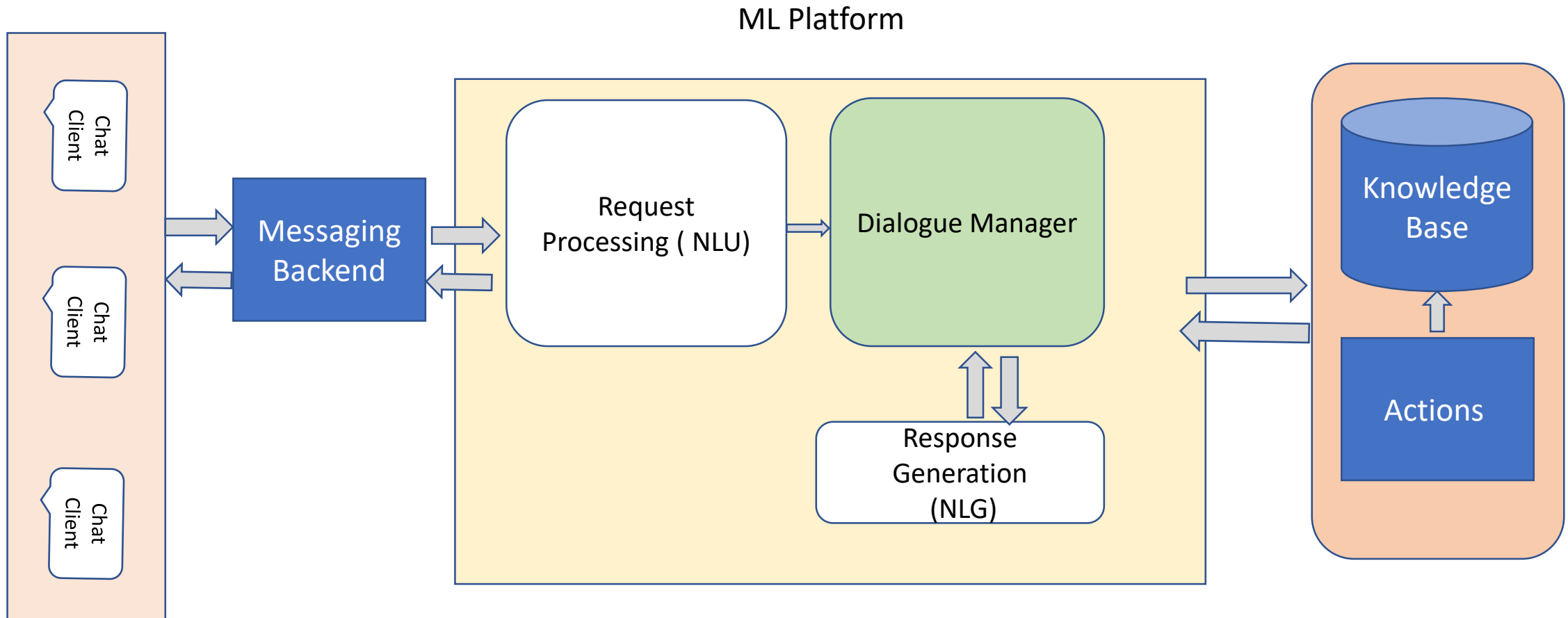
## Rule based

- scripted response to return for a message that matches a specified pattern
- Very limited context and intelligence

## AI based

- Understand context of conversation (Memory)
- Human-like intelligence /Fail gracefully
- NLP, NLG

# Building Blocks



# Intent Classification

Book me a flight from Pune to New Delhi for next week

Need to travel to Delhi this weekend

Planning a trip to Capital city and back this month end

Intent Identified

"Flight Booking"

# Intent Classification

- Intent Classifier is typically a supervised "Text Classification" Problem. Every User message is classified to belong to one of labelled categories (intents)
- Tokenize each word of incoming message
- Use Pre-trained Word embeddings ( BOW Models like GloVe or Word2Vec) to vectorize
- We can also train a model from scratch to generate these embeddings. However, it will require a lot of training data.
- Average individual word embeddings to create embedding for incoming message
- Train a multi-class SVC (Support Vector Classifier or LogReg) to predict probability
- Semantic similar message will appear closer (cosine similarity)

# Entity Recognition

Need to travel to **Delhi** **GPE** **tomorrow** **DATE** morning for Sales meet.

- Extract named entities from the information and categorize them e.g. LOC, DATE, PEOPLE etc
- An entity can be system type (like DATE, email address etc) or application type ( Name of Org)
- Apply NLP Processing pipeline for entity extraction ( Tokenize, POS Tag, Chunking, Entity Extraction)
- A Conditional Random Field ( CRF) model is used to train and detect relevant entities in each query.
- Annotate training data with IOB (or BILUO) and POS Tags for feature generation
- Gazetteer and linguistic features can be used for CRF model

# Dialogue Management

- To manage the conversational aspect (story building)
- uses pattern-based rules to determine the dialogue state for each incoming request
- implements handlers which execute business logic and return a natural language response to the user
- **Slots** – Minimum Input data required to fulfill the goal. **E.g. Flight booking intent requires Source, destination, date and passenger details**
- **Tracker** – Keep track of current state of conversation ( slots, event history)
- **Policy** – What action to take next based on tracker & inputs
  - Fallback Policy
  - Max History to use
- LSTM based RNN models are used for predicting next action which bot will take



# Dialogue Handling - Rules

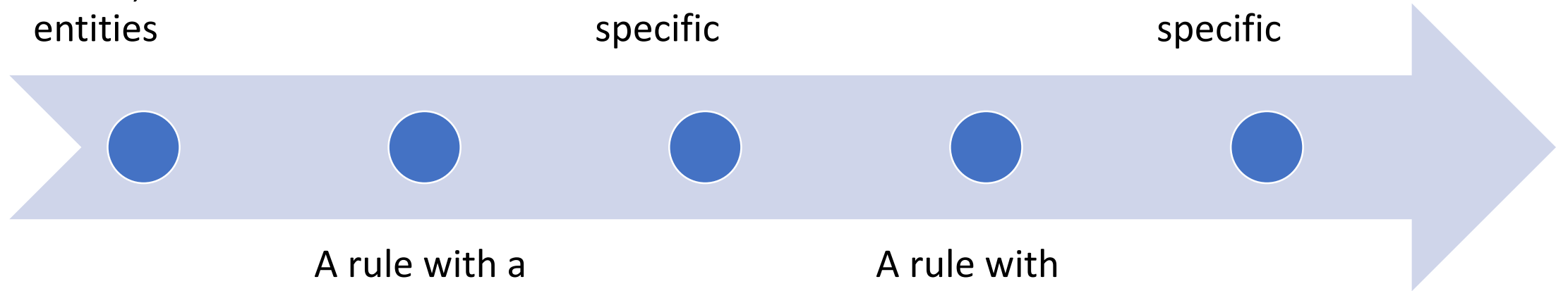
without  
domain,  
intent, or  
entities

A rule with  
an intent is  
more  
specific

A rule  
with *the*  
*most* entities  
is the most  
specific

A rule with a  
domain has  
some  
specificity

A rule with  
entities is  
still more  
specific



# Accuracy of your chatbot

## **Business KPI**

- Goal completion rate - User was able to perform intended action successfully
- Repeat-User – Repeat user after interacting with bot
- Fallback-Rate – Bot decided to transfer conversation to human agent

## **ML Model KPI**

- Confusion Matrix - Precision / Recall / F1 Score

# Challenges

- Handling OOV Words ( Out of Vocabulary Words)
  - Train model with custom domain corpus
  - Merge custom corpus with existing word vectors
- Multiple Intent within same message
  - Message - "Hey! Looks a lovely morning, how is weather in London today"
  - Intents Identified – Greet and Weather Forecast
- Adjusting threshold value
  - A message is identified with multiple intents and relative probability
- Skewness
  - Training data is biased towards a class i.e. more training data is available for one class.
- Limited training data
  - We can't pre-feed all possible routes conversation can take place, to train
- Deal with Open-ended questions
- Adapting to user mood

Thank You.

Questions?

# NLP Internals

Multiple libraries and services are available for common NLP tasks - Spacy/NLTK/Text Blob

Spacy – Popular python library for NLP tasks.

## Linguistic Features

- Word Tokenizer
- POS Tagging
- Lemmatization / Stemming
- Dependency Parsing
- Topic Modeling

For processing textual data, words need to be represented as vectors. Once represented as vectors, we can establish relationship or similarity between words which can be used to train the model.

Word embeddings are just vectors of 300 or 400 floats that represent different words, but a pretrained language model not only has those, but has also been trained to get a representation of full sentences and documents.

Multiple techniques are used

- Count Vectorizer
- TF-IDF ( Term frequency / Inverse Document Frequency)
- CBOW ( Common Bag of Words)
- Skip-gram

\* Cosine-Similarity

GloVE / Word2Vec – Pre-trained word embeddings of most common words made available by Google. Can be load into Spacy

NLU Pipeline

