# About Myself

- **Current Role**: Senior Applied Scientist @ Jupiter
- **Experience**: Conversational AI, Information Retrieval, Information Extraction, Recommendation Engines
- **Education**: MSc. Artificial Intelligence, King's College London, MSc. Tech. Information Systems, BITS Pilani
- **Beyond work**: Travel, music, tea enthusiast!

# Agenda

1. Search engines vs semantic search engines
2. Science behind semantic search engines
3. Framework for building a semantic search engine
4. Why is it such a hard problem in the real world?
5. How to improve such systems?

# Haven't used a search engine?

# What are search engines?

# Search Engine

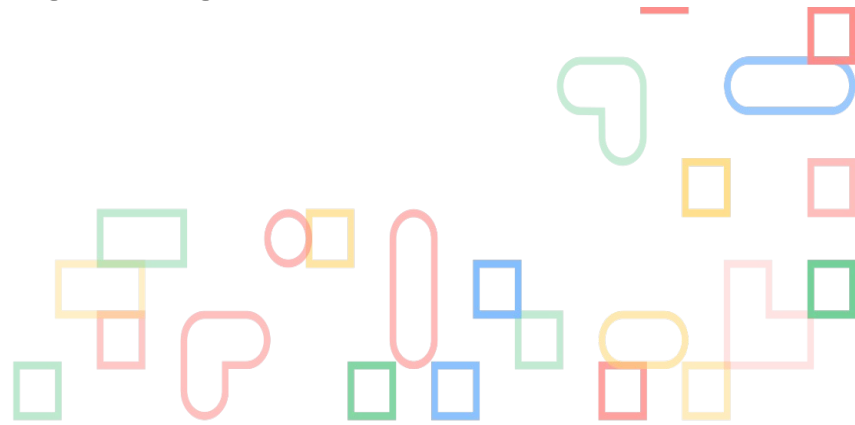Helps people find the information they are looking for online using keywords or phrases.

# What is semantic search?

# Semantic Search

- Understands the intent of the user / searcher via the contextual meaning of search text to generate more relevant results.

- Maximises the possibility of users getting the best search experience possible.

# Does Google use semantic search?

# Example



Keyword Search

Semantic Search

# Another example



Keyword Search

Semantic Search

Google Developer Groups
Cloud • Kolkata

# Underlying problem?

# Information Retrieval



A systematic approach to information retrieval, Lalmas et al. 2001

# AI systems = Code + Data

# Data for ad hoc IR

1. Corpus of search queries
2. Corpus of candidate documents
3. Ground truth (explicit or implicit)

# Solution Architecture

# Semantic Search Architecture

**Offline** steps:

1.  Document representation
2.  Document indexing

**Online** steps:

1.  Query representation
2.  Candidate retrieval
3.  Candidate ranking

# Text representation

- Sparse features: 1 hot representations of characters (TFIDF), character n-grams

- Dense representations: Term or pre-trained embeddings

# How to choose a representation ?

# How to choose the right representation?

- Depends on size, property and domain of the data available
- DNNs for IR can learn text representations in situ, or use pre-trained embeddings.
- For pre-trained embeddings, make sure the representations are suitable for the task.
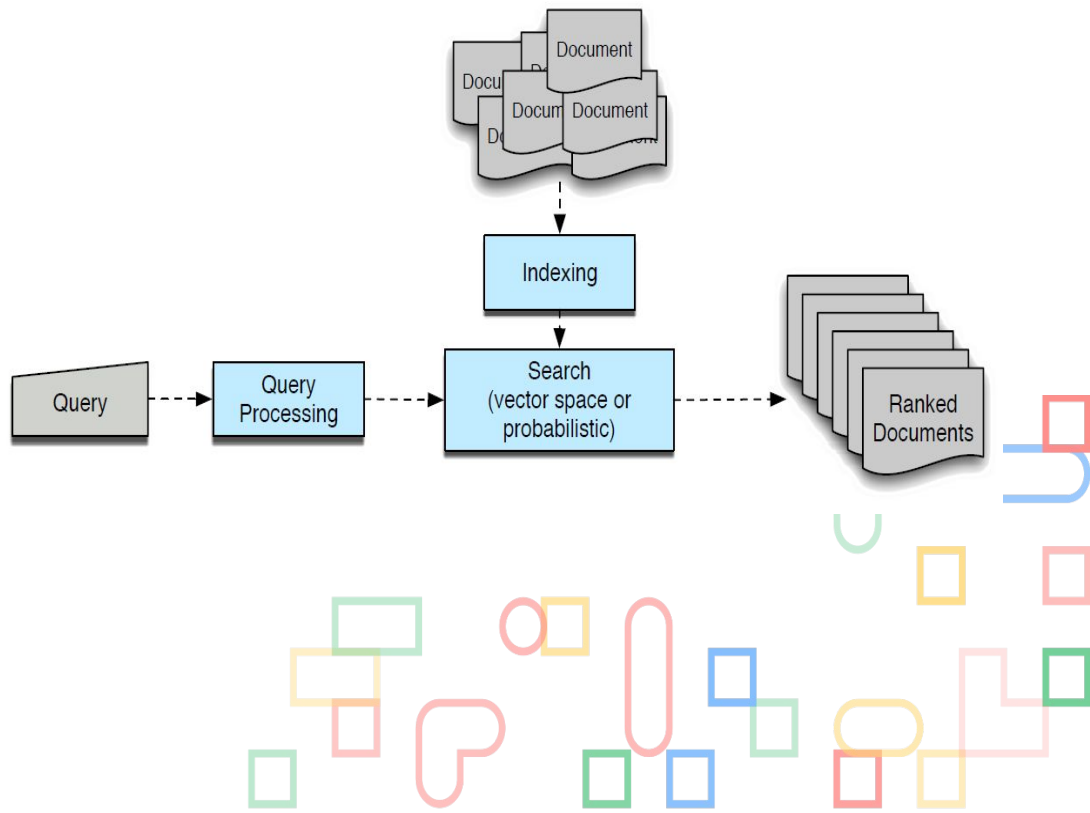
Embedding based models often make different errors than exact matching models and the combination of two may be more effective.

# Document Indexing

Linking or tagging documents with certain attributes which makes discovering them easier.

# Candidate retrieval

- Retrieve matching candidates to the query on the basis of some measure of relevance.

- Rank the retrieved documents by the relevance score.



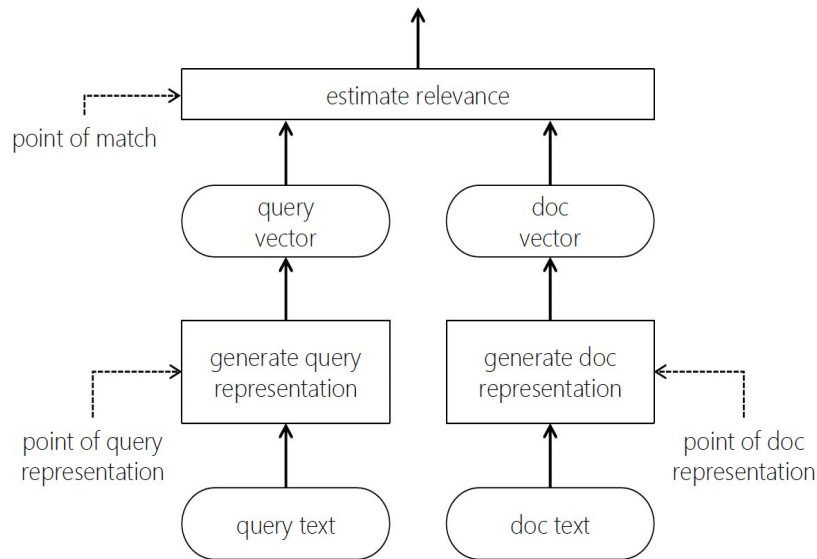Bhaskar Mitra and Nick Craswell (2018), "An Introduction to Neural Information Retrieval",

Google Developer Groups
Cloud • Kolkata

# No one size fits all!

# Retrieval Approaches

**Unsupervised**
- No labelled data available
- Use existing retrieval model or query-document similarity metric (edit distance)

**Semi-supervised**
- Small amount of labelled data available
- Train a retrieval model with few parameters

**Fully Supervised**
- Large amount of labelled data available
- Train a model that optimizes directly for the target task

# Ranking

# Ranking Approaches

- No labelled data available ⇒ Unsupervised

- Labelled data (ranked ground truth) available
  ⇒ Supervised Learning to rank (LTR)

Google Developer Groups
Cloud • Kolkata

# You can't improve what you can't measure!

# Metrics

**Offline**
- MAP (Mean Average Precision)
- MRR (Mean Reciprocal Rank)
- NDCG (Normalized Discounted Cumulative Gain)

**Online**
- A/B testing
- Click through rate
- View rate
- Any other business metric

# Importance of feedback loops

# Why is search so difficult?
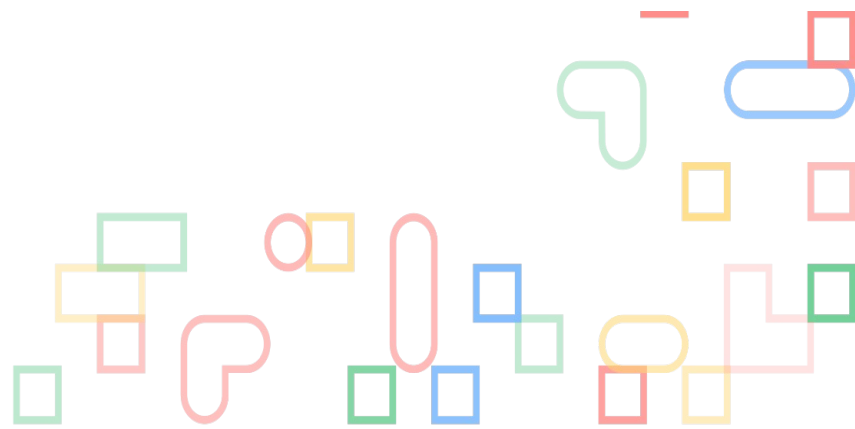
# Unstructured Data

- Data drift
- Spelling errors
- Domain adaptation
- Multilingual search queries
- Lack of data (labelled / unlabelled)
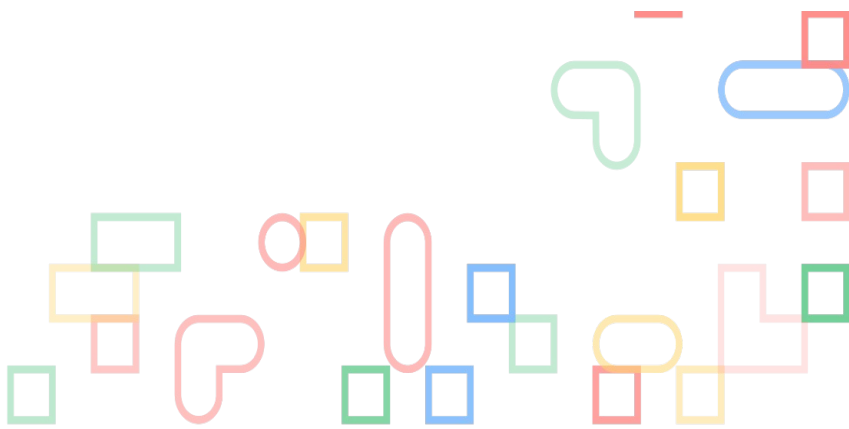- Out-of-scope queries
- No match found

# Tips to improve search

- **Monitor data drift** and use **active learning** feedback loops
- **Autocomplete feature** to prevent spelling mistakes
- Use **transfer learning** to fine tune large pre-trained models with unlabelled domain specific data
- **Analyse search query data** to gauge the need of a multilingual embedding model (Use 80/20 principle!)
- Explore **data augmentation/synthetic data** creation techniques
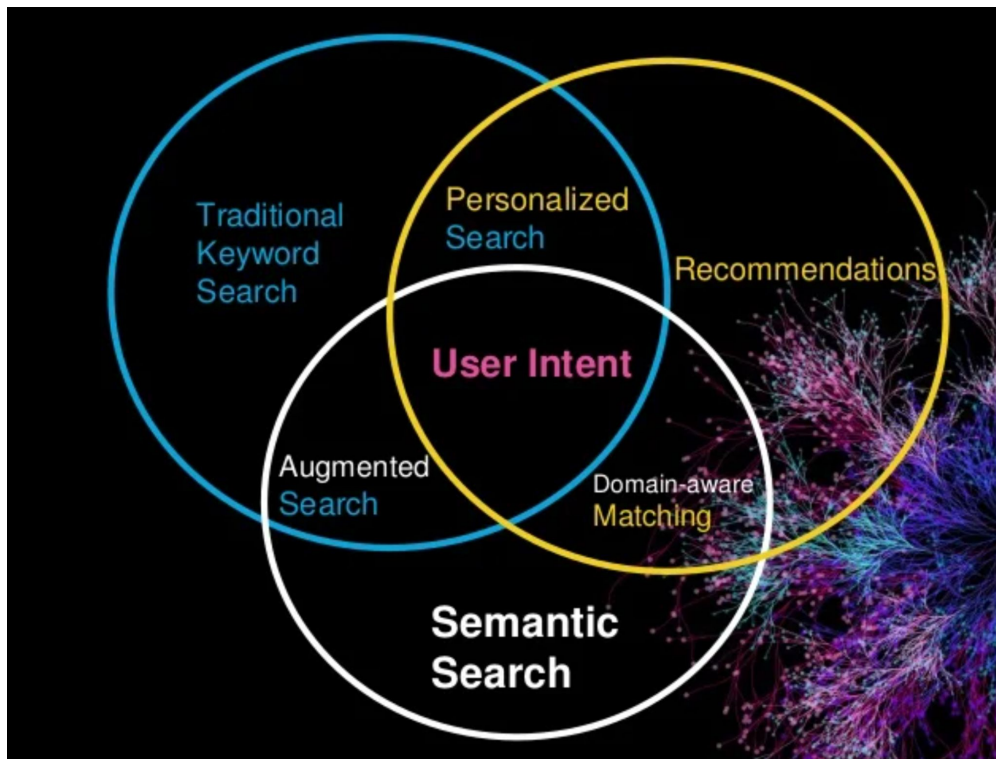- Handle **no match found/out-of-scope** queries gracefully

# Tips to improve search

- **Diagnose/improve content** when users aren't clicking on search results
- **Feature sought after content prominently** for everyone to see (most searched content)
- Analyze the page in the app the user was visiting just before they started searching, did the page provide the information a user searched for, **poor UX design**?

# Beyond Semantic Search

# Thank you!

## Questions?

@shweta_bhatt8