

ElastiK Nearest Neighbors

An Elasticsearch Plugin to Simplify
Online K-Nearest-Neighbors Search

Alex Klibisz - Insight Data Engineering - Boston

Twitter Image Similarity Search

Searched 5703728 images in 171 milliseconds

Query Image



Original Tweet

Nearest Neighbors



Application: Image Similarity Search Engine

Any Online KNN Application

E.g. similar item recommendation for e-commerce

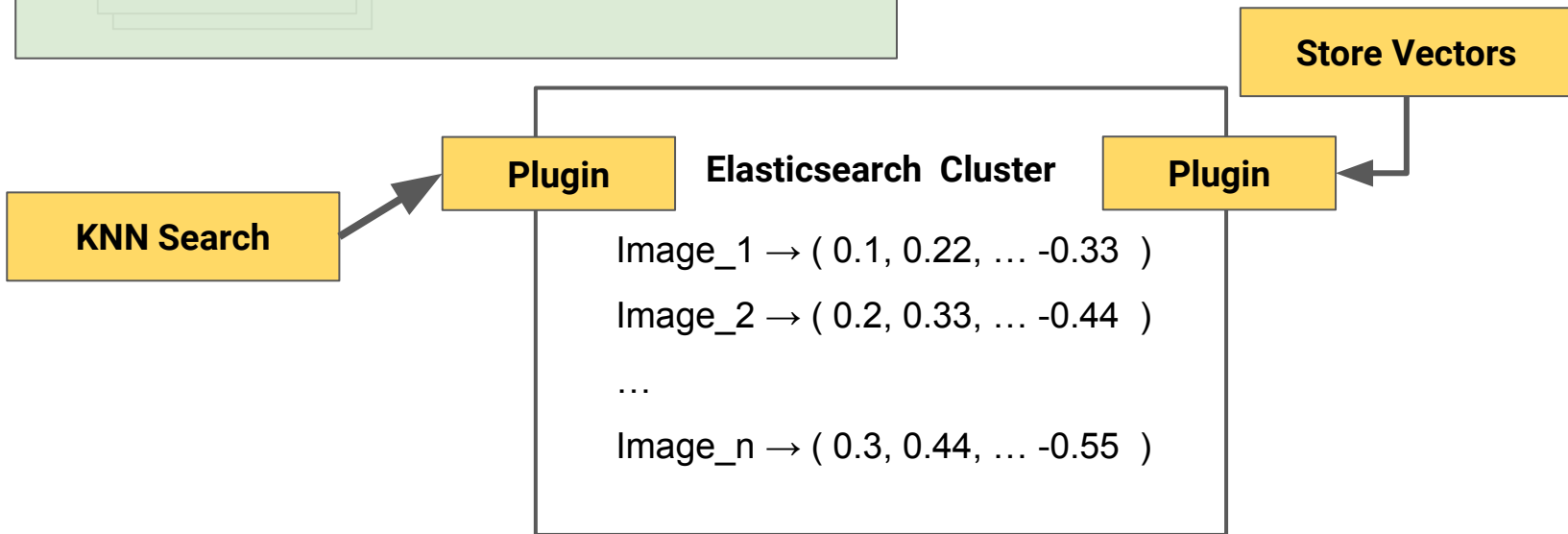
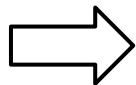
Feature vectors preserve similarity

Image_1 \rightarrow (0.1, 0.22, ... -0.33)

Image_2 \rightarrow (0.2, 0.33, ... -0.44)

...

Image_n \rightarrow (0.3, 0.44, ... -0.55)



KNN on Elasticsearch - Why is it useful?

KNN in **offline** setting → standard batch infrastructure:

- Static corpus of vectors
- Batch job computes and caches neighbors

KNN in **online** setting → complicated infrastructure:

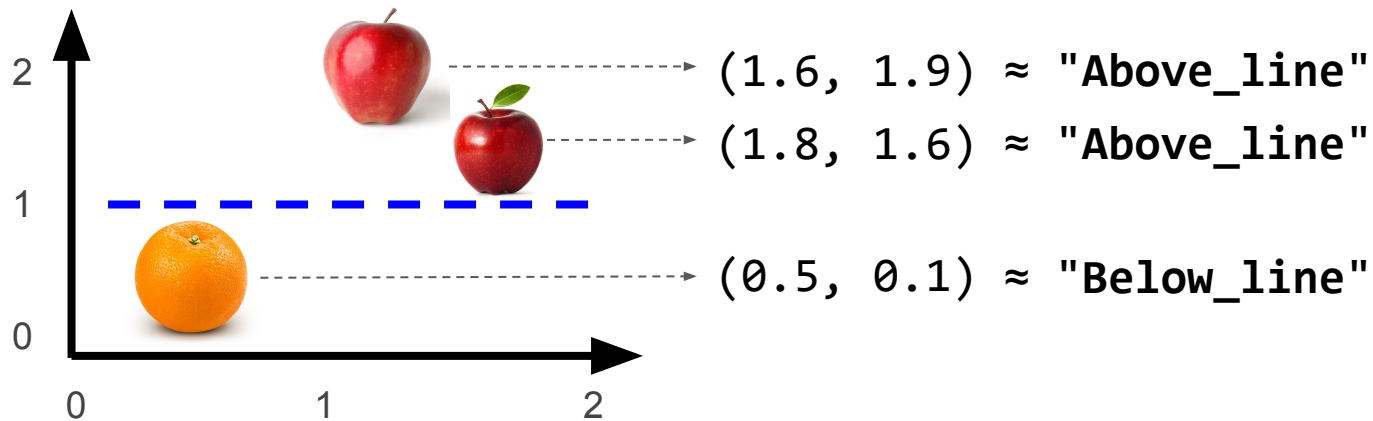
- Store and randomly access millions of vectors
- Constantly add more items to corpus
- Scale horizontally for many concurrent searches

Put the nearest neighbors search on a proven infrastructure: Elasticsearch

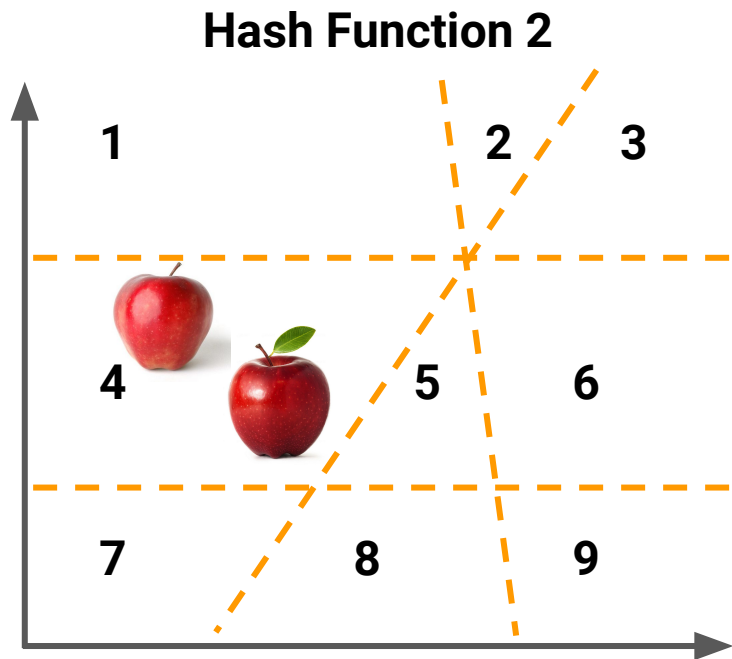
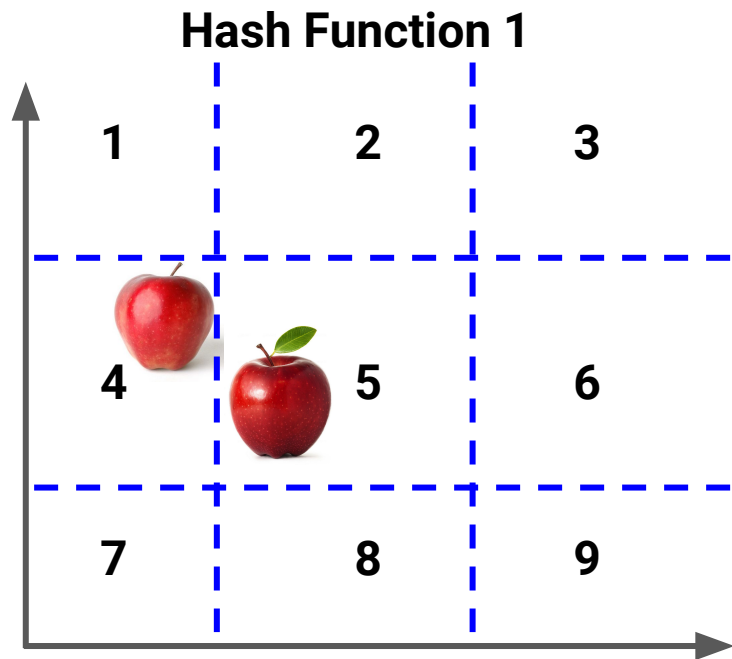
Indexing and searching floating-point vectors?

Locality Sensitive Hashing

- Represent vectors as discrete tokens while preserving similarity



Mapping LSH to Elasticsearch



More hash functions → More likely similar vectors share regions

How to store, query for vectors that share regions?

Mapping LSH to Elasticsearch

Query vector

```
"hashes": {  
  "h1": 8  
  "h2": 31  
  "h3": 44  
}
```

$\text{Hash}_1(\text{vec}) = 8$
 $\text{Hash}_2(\text{vec}) = 31$
 $\text{Hash}_3(\text{vec}) = 44$

BooleanQuery

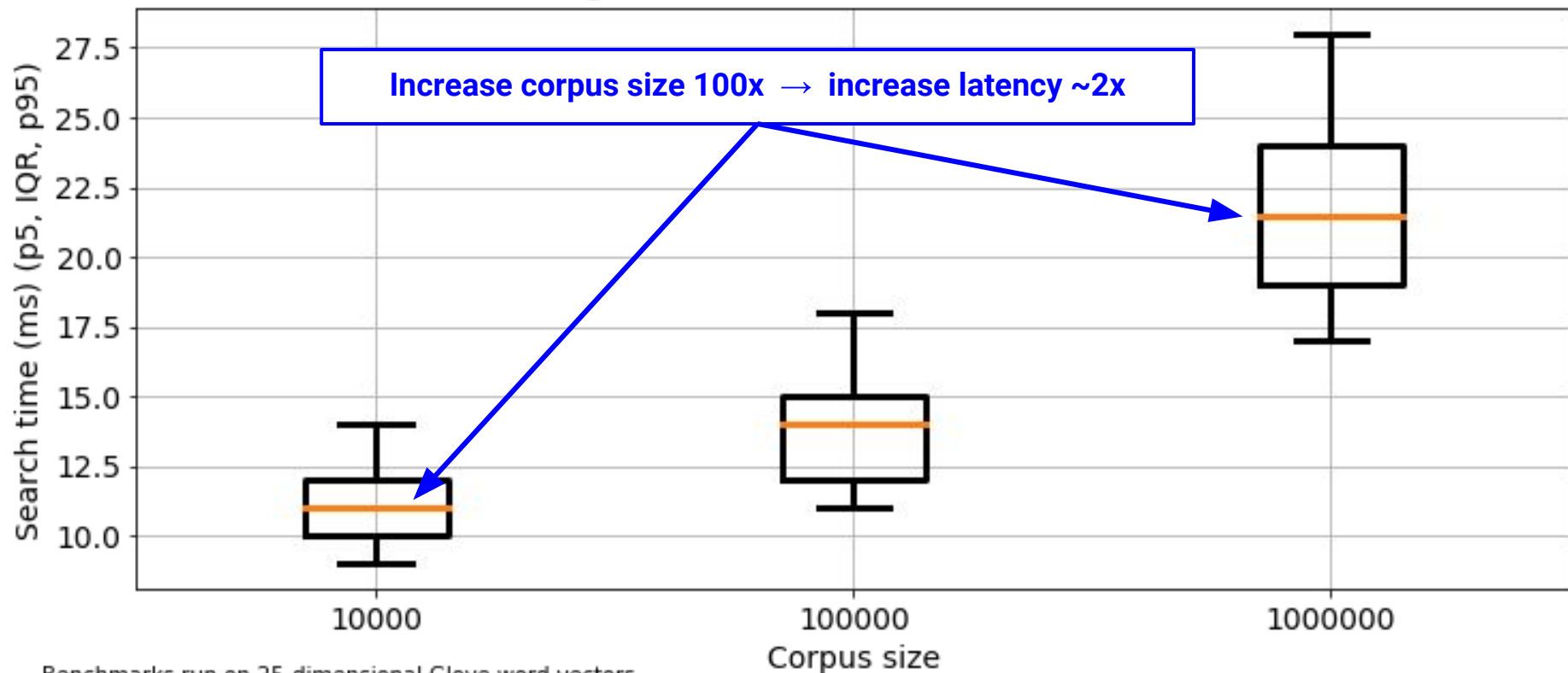
- Return vectors with most matching key-value pairs
- Search time scales sub-linearly with respect to corpus size
 - Critical property for online setting

Candidate vector

```
"hashes": {  
  "h1": 8  
  "h2": 31  
  "h3": 33  
}  
  
"_similarity": 2/3
```

Search time as a function of corpus size

config=(t=10, b=8, $k_1 = 500$, $k_2 = 10$)



Benchmarks run on 25-dimensional Glove word vectors
(<http://nlp.stanford.edu/data/glove.twitter.27B.zip>)

Boxplots represent 1000 sample searches

More details available

Implementation

Performance (speed and recall)

Image processing pipeline

... saved for Q and A, see supplementary slides

Alex Klibisz

Computer science, University of Tennessee

Enjoy productionizing ML Systems

Enjoy tennis, mountain-biking, traveling, and podcasts



Kayaking in Seward Alaska, May '17

Supplementary Material

Nearest Neighbors Plugin

- Implementation
- Parameters
- Benchmark - search latency
- Benchmark - recall
- Distributed performance
- Time complexity
- Locality sensitive hashing details
- Boolean query details

Full Pipeline

Image processing pipeline

ES-Aknn Implementation

Java with Gradle build system

Operates as a middleware

- All logic runs on Elasticsearch nodes
- Register endpoint handlers on start
- Parse HTTP request parameters and JSON requests when endpoints are hit
- Query and store documents using ElasticSearch Java API
- Construct and return JSON responses

ES-Aknn Parameters

`nb_tables`

- Number of hash functions applied to each vector
- Increasing this increases recall, decreases speed

`nb_bits_per_table`

- Each hash function partitions the space into $2^{\text{nb_bits_per_table}}$ buckets
- Changing this affects recall, shouldn't affect speed

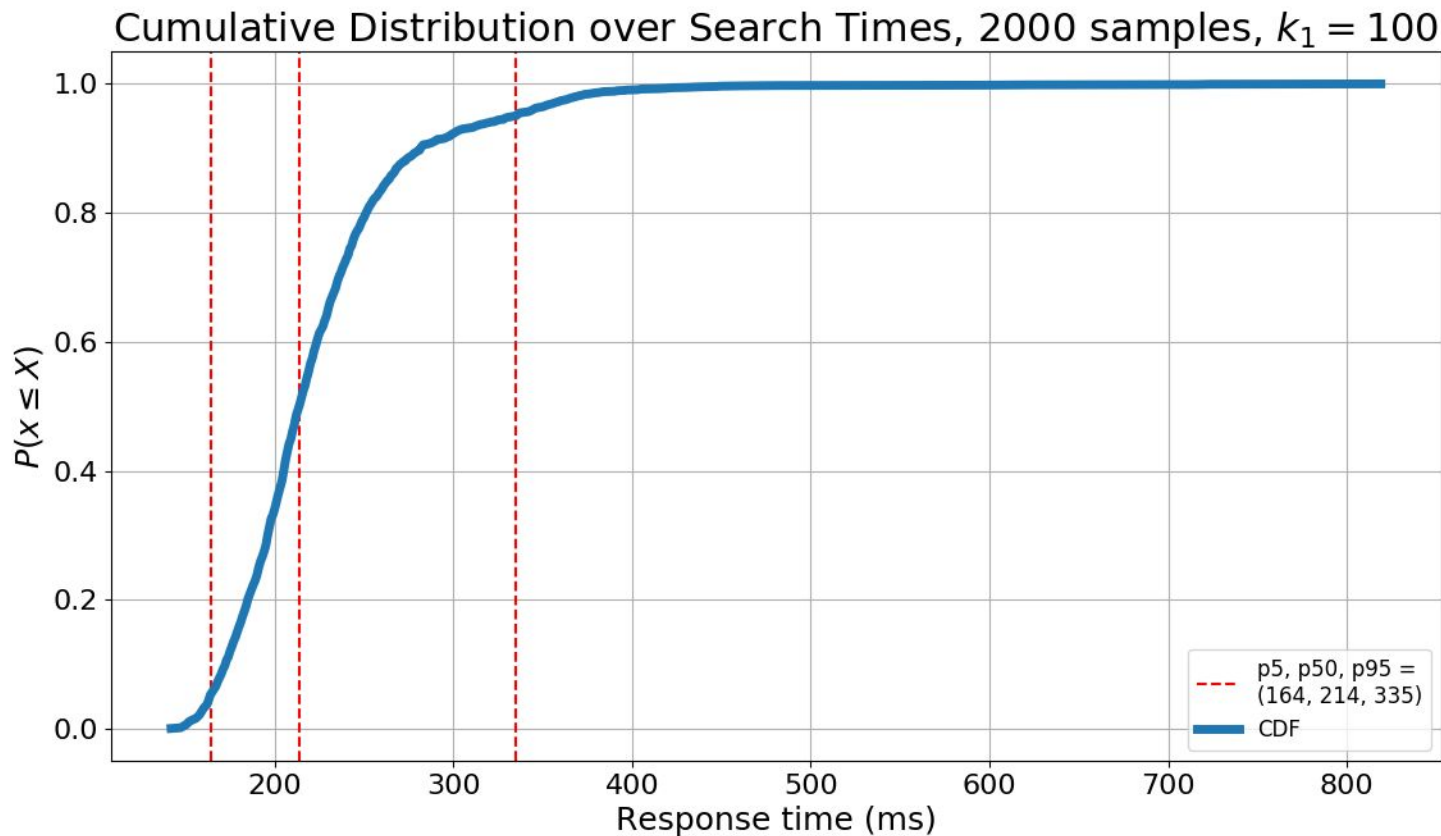
`k1`

- Number of approximate neighbors considered for exact computation
- Increasing this increases recall, decreases speed

`k2`

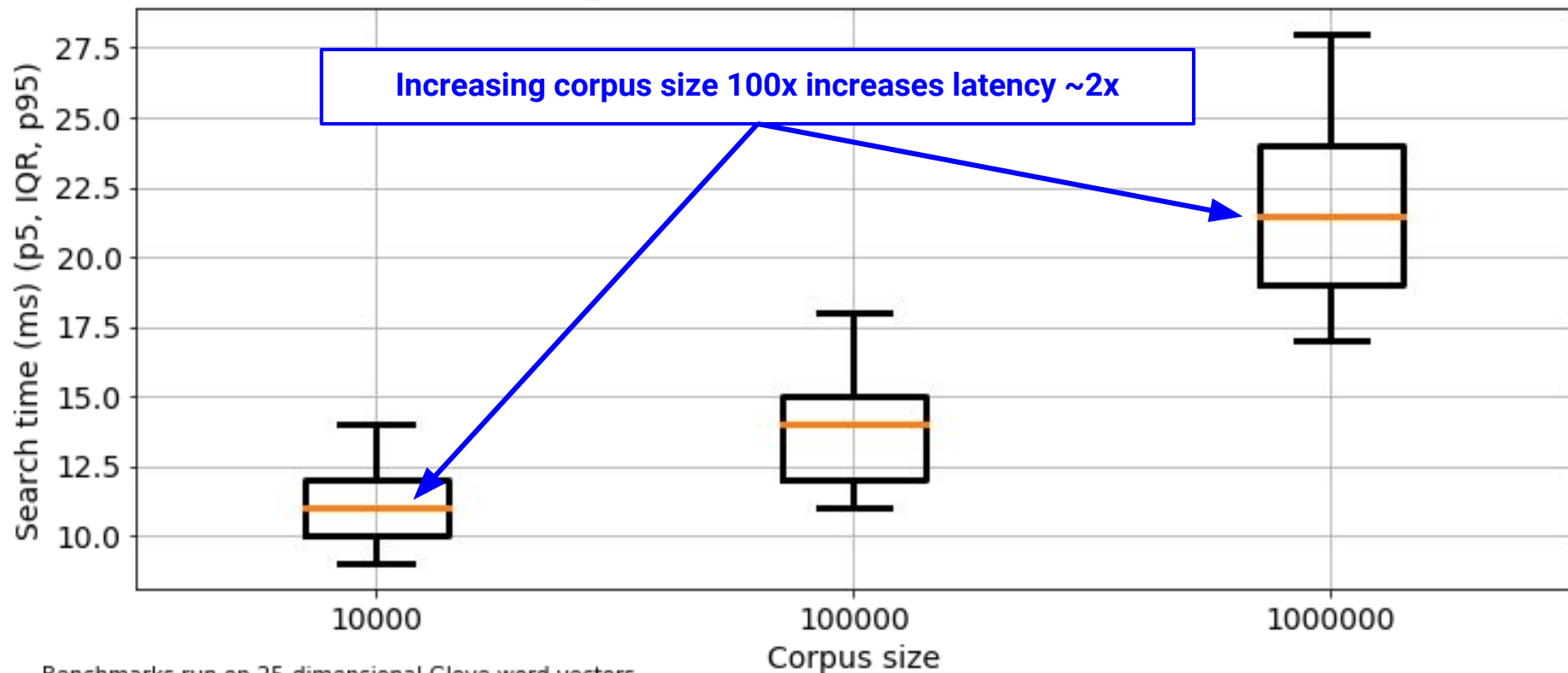
- Number of exact neighbors returned
- Changing this has generally negligible effect on speed

Image Similarity Search Times



Search time as a function of corpus size

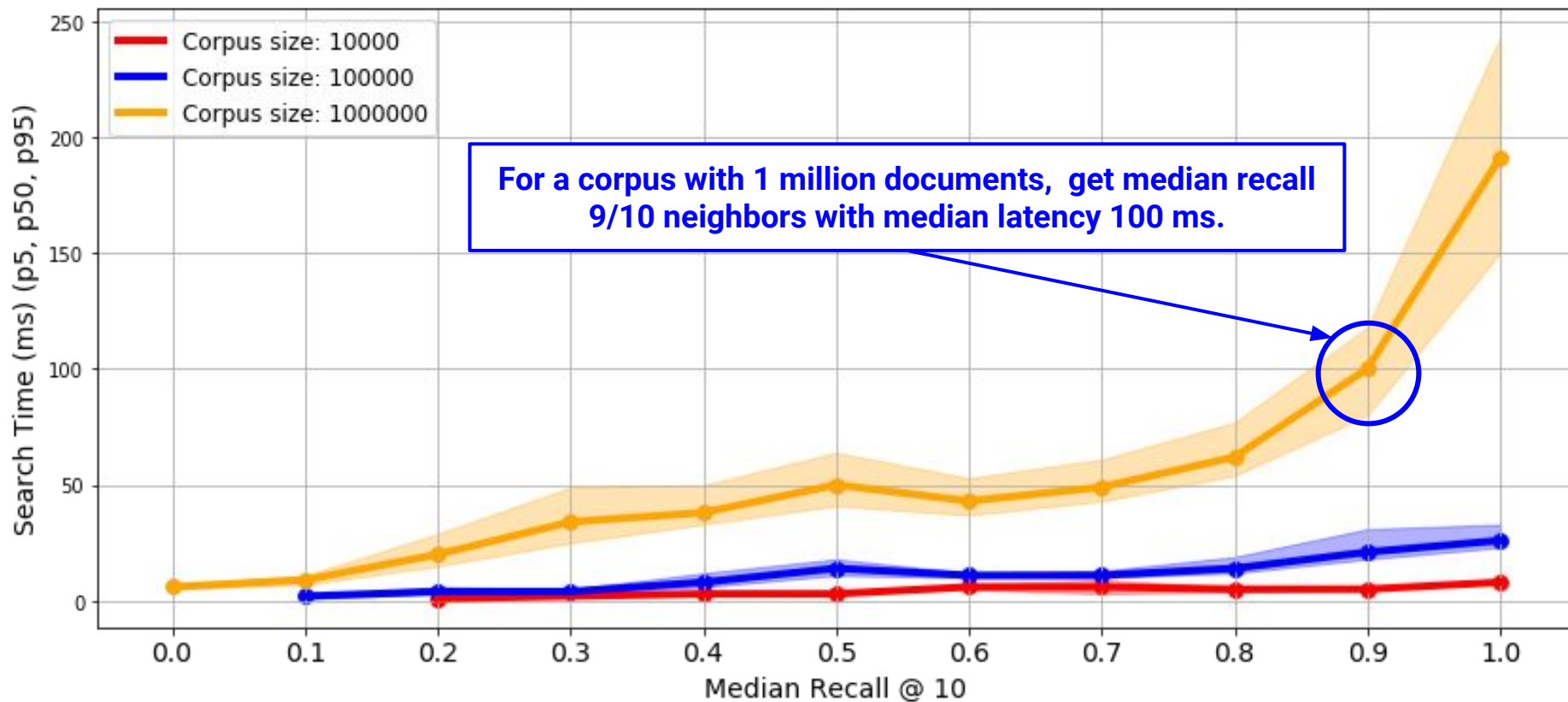
config=(t=10, b=8, $k_1 = 500$, $k_2 = 10$)



Benchmarks run on 25-dimensional Glove word vectors
(<http://nlp.stanford.edu/data/glove.twitter.27B.zip>)

Boxplots represent 1000 sample searches

Search time as a function of recall



For a corpus with 1 million documents, get median recall 9/10 neighbors with median latency 100 ms.

Benchmarks run on 25-dimensional Glove word vectors
(<http://nlp.stanford.edu/data/glove.twitter.27B.zip>)

Each point represents 1000 sample searches

ES-Aknn Distributed Performance

Setup

- 4x M5.xLarge (4-core, 16GB, SSD), 5x ES shards
- 300-dimensional Glove vectors

Indexing new documents

- 1600 new docs / second using round-robin posts containing 10K vectors

Concurrent searches

- 100 concurrent queries, LSH with 32 tables, 32 bits
- Latency for corpus of 1M vectors: mean \approx 300 ms, stdv \approx 80 ms

ES-Aknn Complexity

Search

$O(\text{BooleanQuery}(N, \text{nb_tables}, k1) + (k1 \times \text{nb_dimensions}) + \text{partial-sort}(k1, k2))$

Find $k1$ candidate vectors with matching hashes compute exact similarity find the top $k2$

$\equiv O(\text{BooleanQuery}(N, \text{nb_tables}, k1)) \ll O(N)$

In practice the Lucene BooleanQuery dominates runtime. It's typically logarithmic or better but depends on configurations.

Indexing

$O((\text{nb_tables} \times \text{nb_bits_per_table} \times \text{nb_dimensions}) + \text{Index}(\text{nb_tables}))$

Hash the floating-point vector via dot products

Index as a standard ES document

The time spent on hashing vs. indexing depends on the batch size

LSH Implementation

Build LSH Model from sample of vectors

```
5 def make_lsh_model(nb_tables, nb_bits, nb_dimensions, vector_sample):
6     # vector_sample: np arr w/ shape (2 * nb_tables * nb_tables, nb_dimensions).
7     # normals, midpoints: np arrs w/ shape (nb_bits, nb_dimensions)
8     # thresholds: np arrs w/ shape (nb_bits)
9     # all_normals, all_thresholds: lists w/ one normal, one threshold per table.
10    all_normals, all_thresholds = [], []
11    for i in range(0, len(vector_sample), 2 * nb_bits):
12        vector_sample_a = vector_sample[i:i + nb_bits]
13        vector_sample_b = vector_sample[i + nb_bits: i + 2 * nb_bits]
14        midpoints = (vector_sample_a + vector_sample_b) / 2
15        normals = vector_sample_a - midpoints
16        thresholds = np.zeros(nb_bits)
17        for j in range(nb_bits):
18            thresholds[j] = normals[j].dot(midpoints[j])
19        all_normals.append(normals)
20        all_thresholds.append(thresholds)
21    return all_normals, all_thresholds
```

LSH Implementation

Compute the hashes for a vector

```
24 def get_lsh_hashes(vec, all_normals, all_thresholds):
25     # vec: np arr w/ shape (nb_dimensions, )
26     # hashes: one hash per table.
27     hashes = dict()
28     for normal, thresholds in zip(all_normals, all_thresholds):
29         hsh = 0
30         dot = vec.dot(normal.T) # shape (nb_bits,)
31         for i, (d, t) in enumerate(zip(dot, thresholds)):
32             if d > t:
33                 hsh += i ** 2
34         hashes[len(hashes)] = hsh
35     return hashes
```

LSH Linear Algebra

Given two sample vectors v_1 and v_2 , compute an equidistant hyperplane

- Hyperplane defined by midpoint m , normal vector n

$$m = (m_1, m_2) = (v_1 + v_2) / 2$$

$$n = (n_1, n_2) = v_1 - m$$

Given a new vector v , compute its hash relative a hyperplane defined by m and n

- The hash function should return 1 if the vector is “above” the hyperplane, 0 if the vector is “below” the hyperplane

$$h(v, m, n) = 1[n \cdot v > n \cdot m] \quad (\text{indicator function})$$

- There are $\text{nb_tables} * \text{nb_bits_per_table}$ such hash functions applied, each using an m and n computed from a random sample pair of vectors.
- Applying the hash functions for a single table also yields a binary tree.

ES-Aknn Boolean Query

GET demo/doc/1

```
{
  "_index": "demo",
  "_type": "doc",
  "_id": "1",
  "_source": {
    "h0": 10,
    "h1": 20,
    "h2": 30
  }
}
```

GET demo/doc/_search

```
{
  "query": {
    "bool": {
      "should": [
        {"term": {"h0": 10}},
        {"term": {"h1": 20}},
        ...
      ]
    }
  }
}
```

ES-Aknn Applications

General

- K-nearest-neighbors vector storage
- Low-latency K-nearest-neighbors search with growing corpus

Specific

- Image similarity search (a.k.a reverse image search)
- Audio similarity search
- Recommendation engines (especially when frequently adding/updating user/item vectors)

Full Pipeline

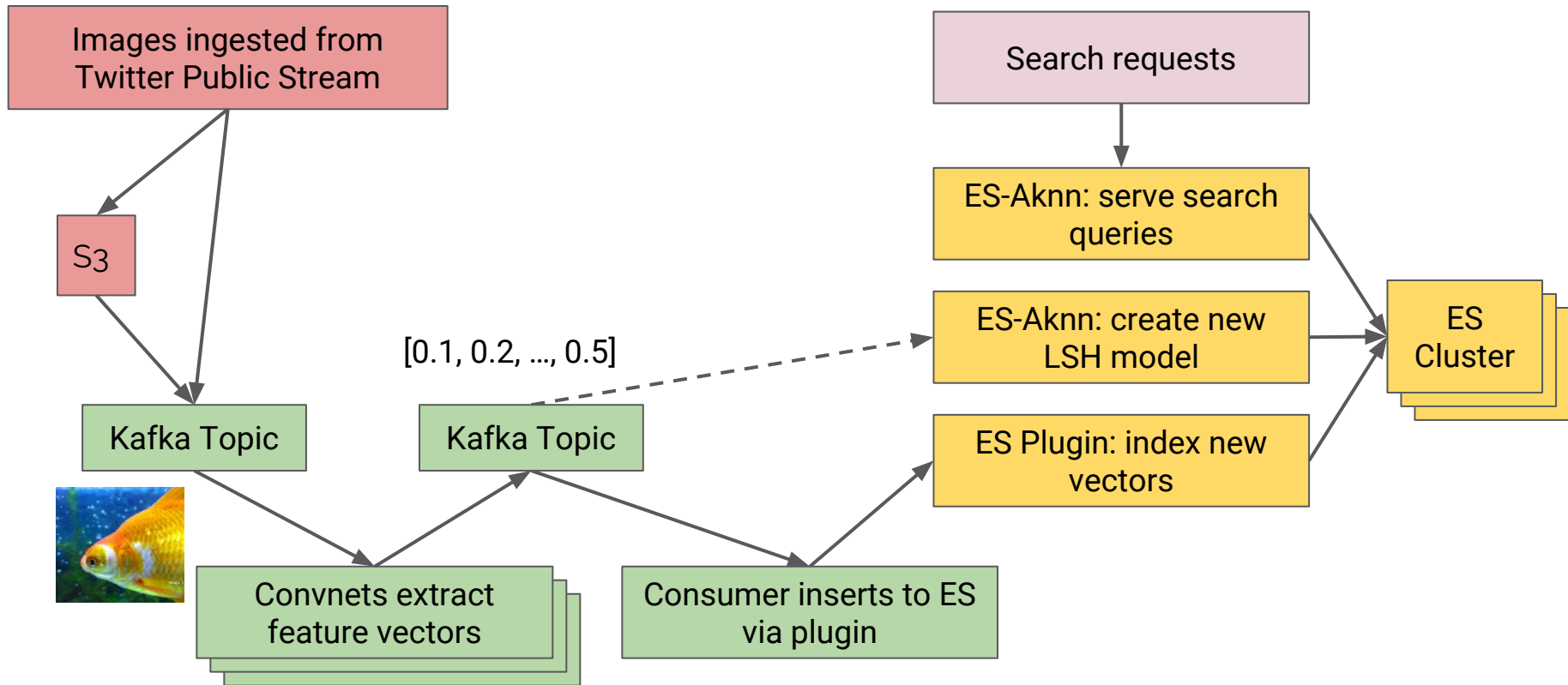


Image Processing Pipeline

Implementation

- Kafka as a distributed job queue with 10 EC2 instance workers
- Python Kafka consumers computing feature vectors
- Keras (MobileNet) pre-trained convnet - use conv_preds layer for output
- Consumers pull images from S3 (faster than funneling through Kafka)
- Parallelized S3 requests (ThreadPoolExecutor), image preprocessing (Multiprocessing Pool)

Performance

- 40 images / node / second with K80 GPU (\$0.3/hr spot instance)
- 33 images / node / second with 36-core CPU (\$0.6/hr spot instance)

Helpful Resources

[LSH lectures by Victor Lavrenko](#)

[ES plugin cookiecutter template](#)

[ES Ingest-OpenNLP plugin](#)

[Presentation about ANNOY by Erik Bernhardsson](#)