

Vector Databases and ANN algorithm





Keyword search

vs Semantic search

A search bar with the word "smartphones" entered and a magnifying glass icon on the right.

exact term matching

synonyms and other
similar/related terms

“Top **smartphones** on the market...”
“**Smartphone** industry struggles...”

“Using a **cell phone** in Spain...”
“**Mobile device** usage on the rise...”

Images

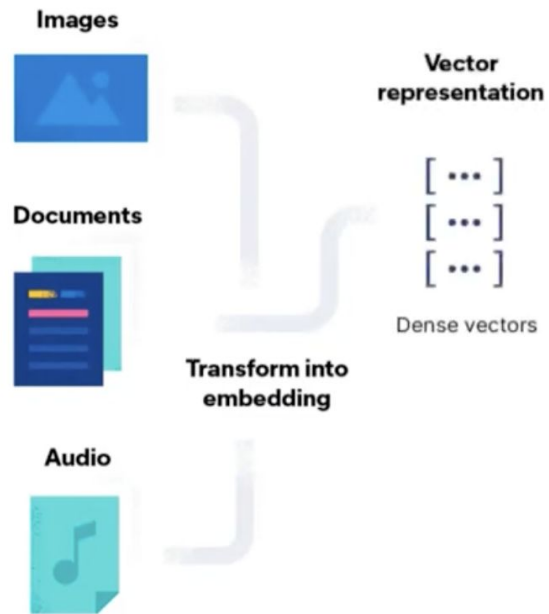


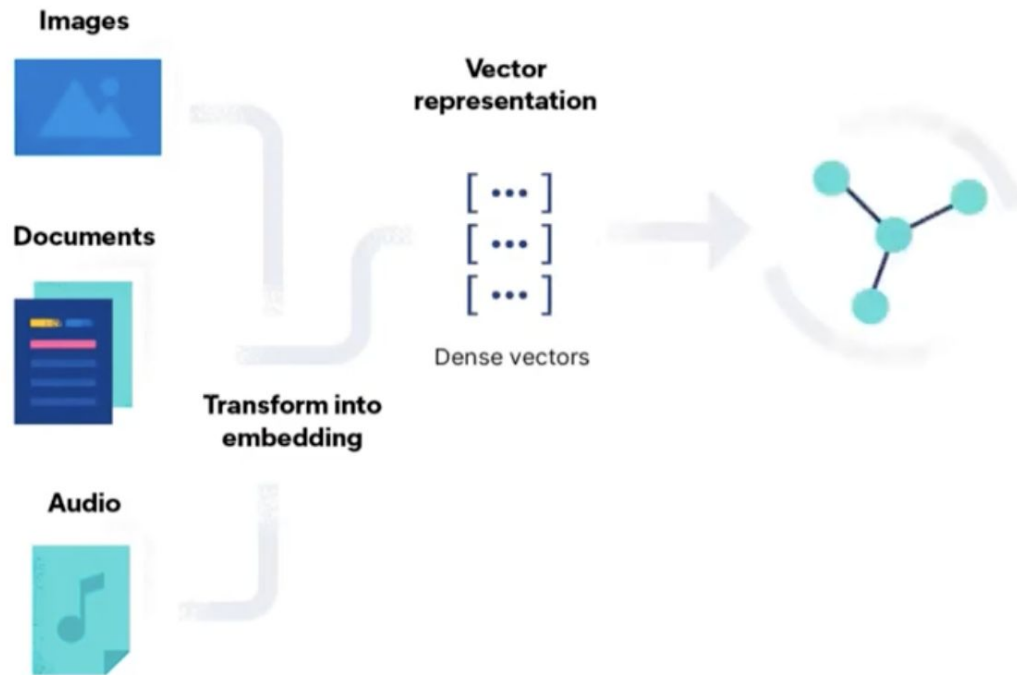
Documents

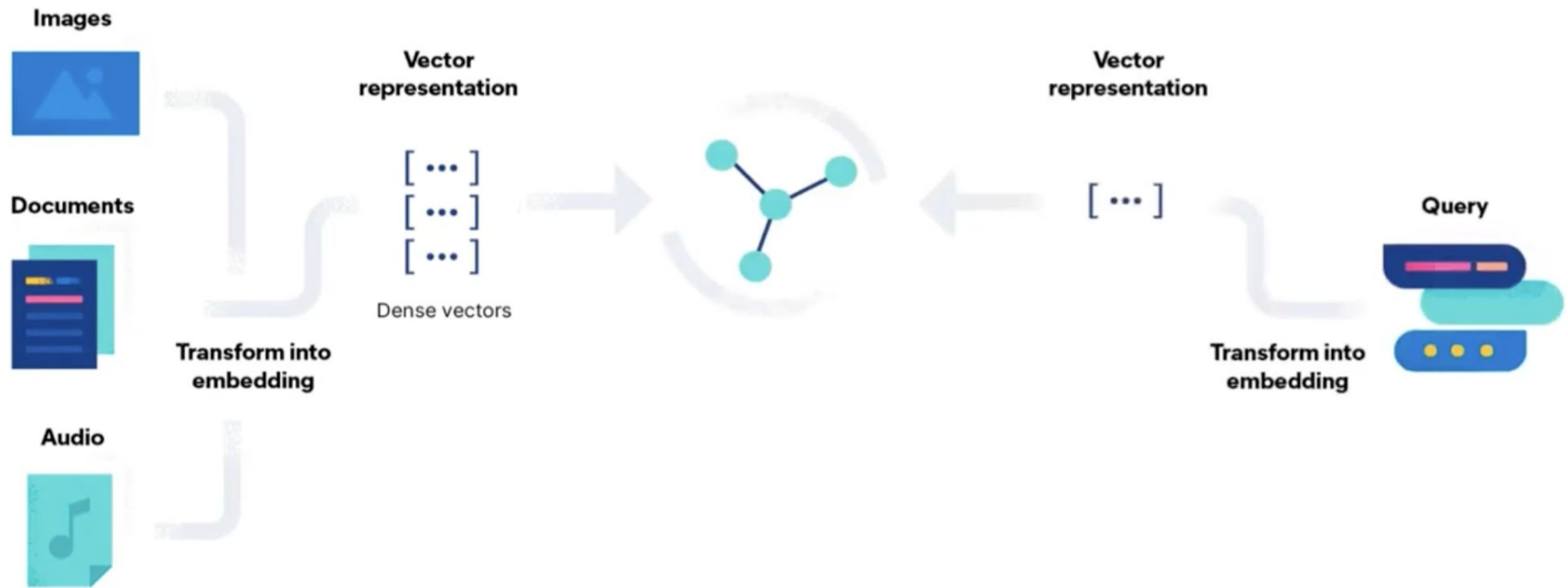


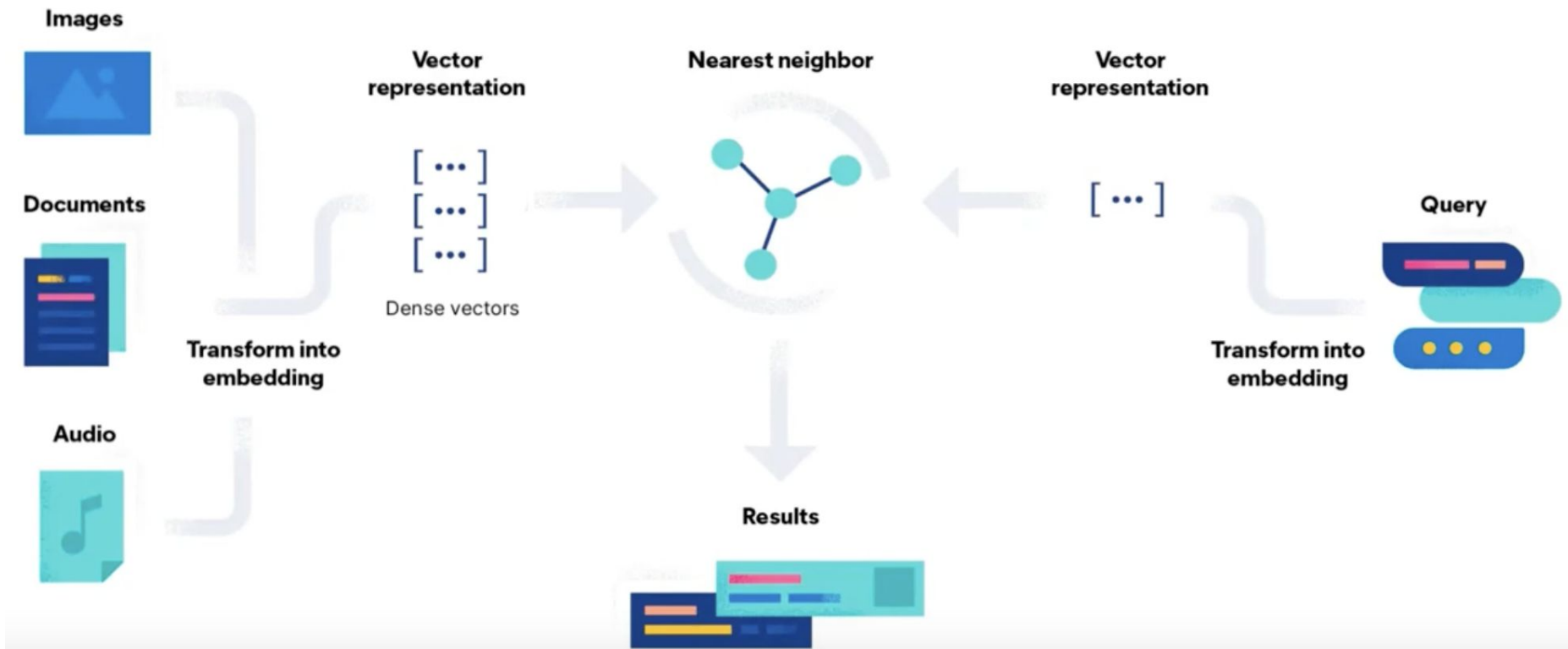
Audio

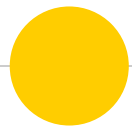












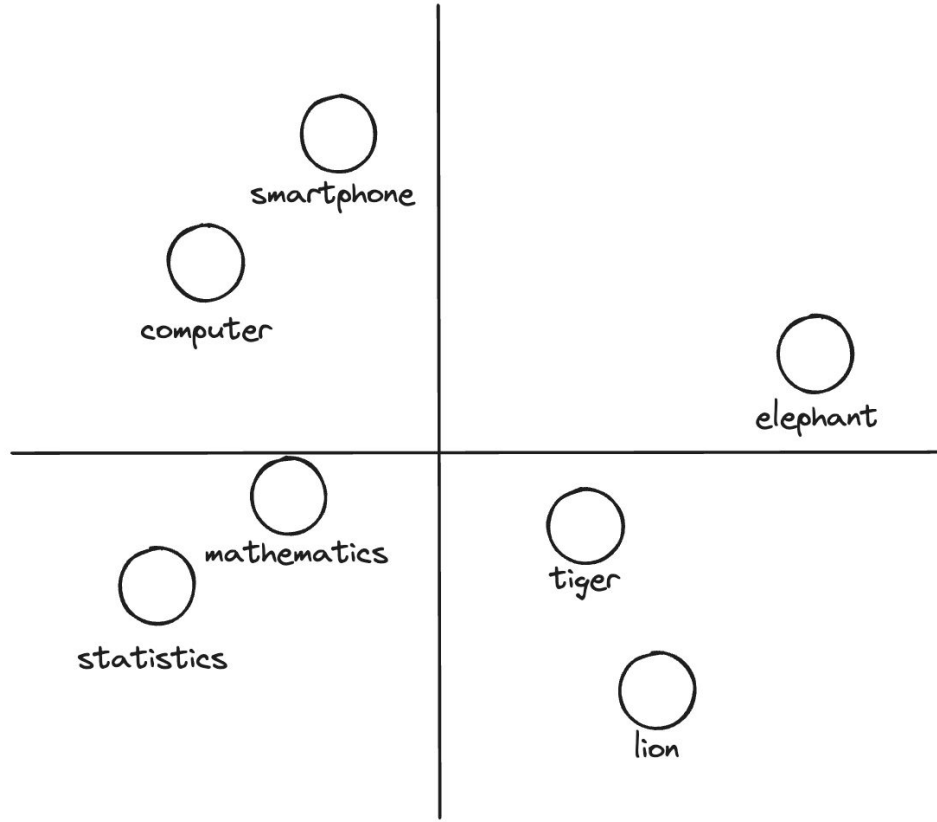
Semantic search is Vector search

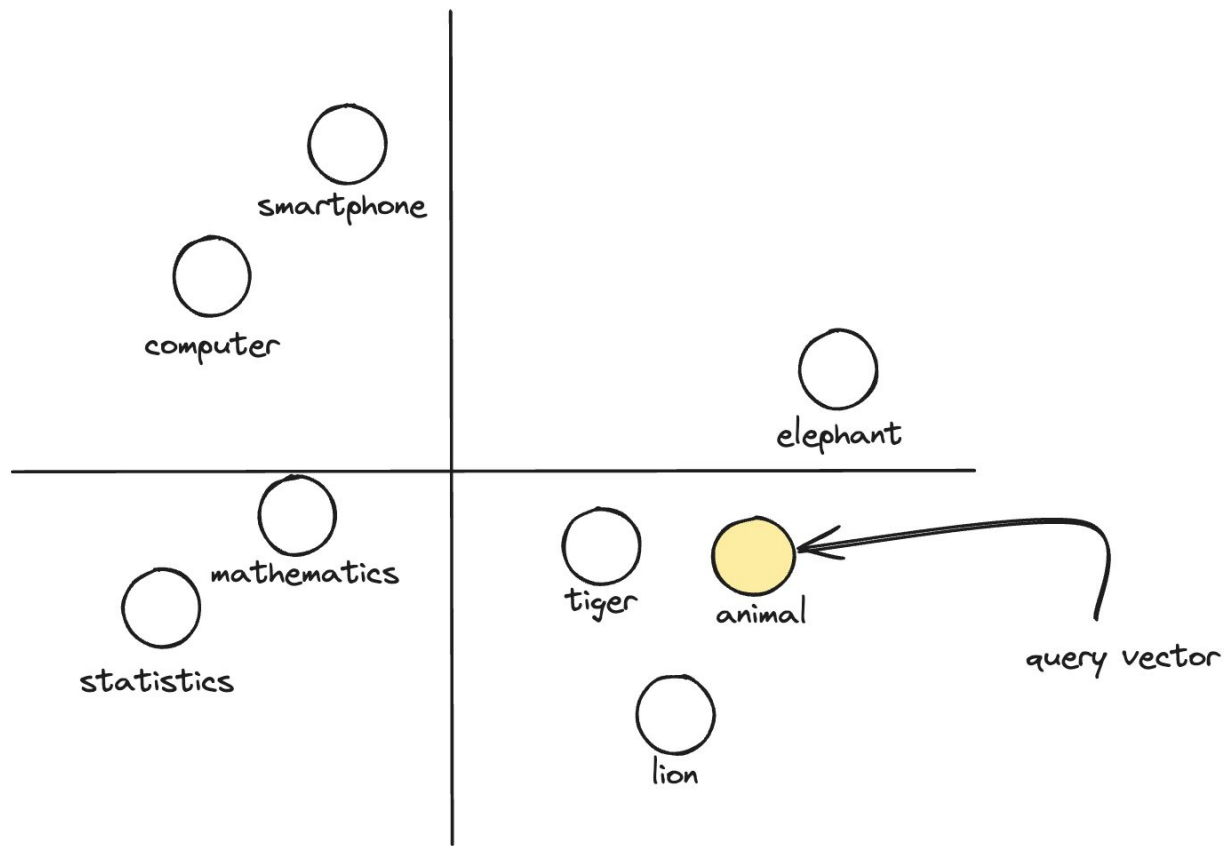


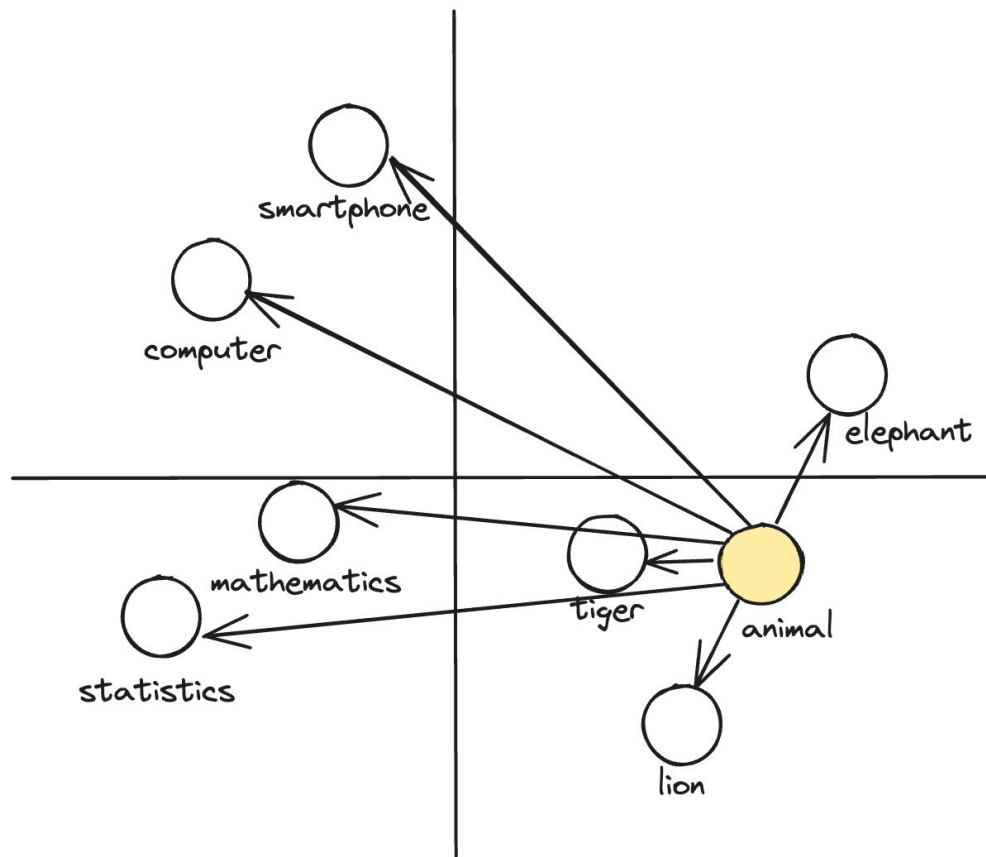
K Nearest Neighbour

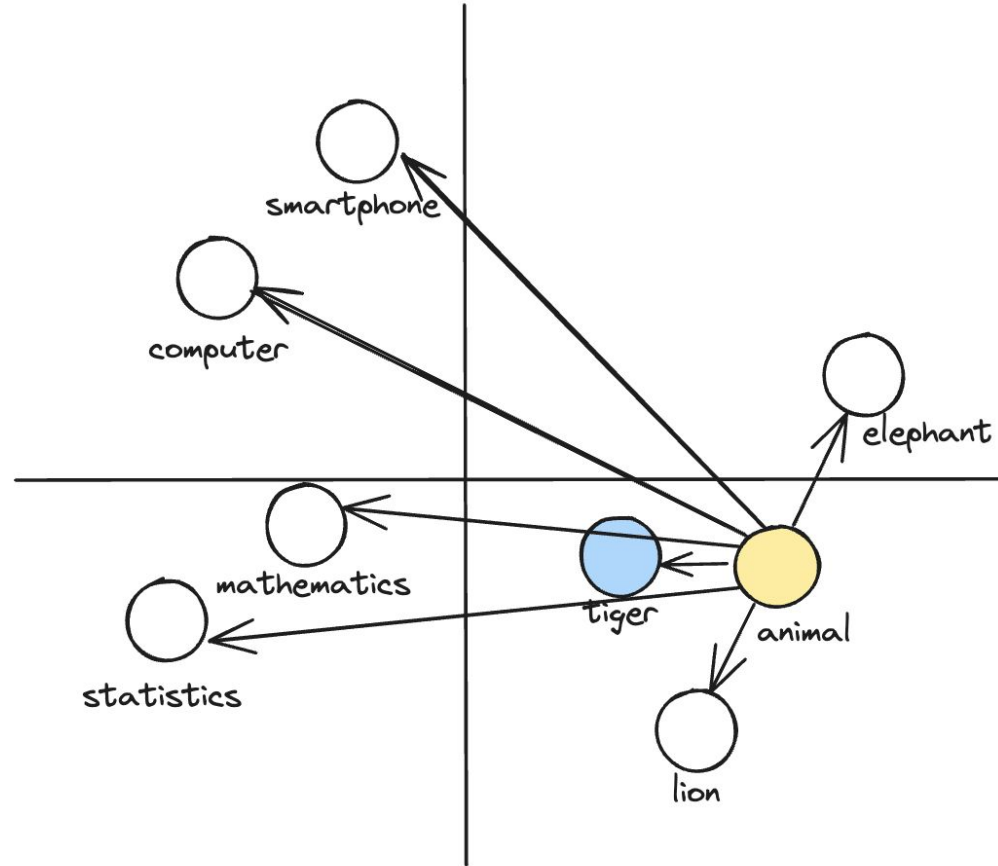
Exhaustive, brute-force approach

- Measure distance from query and all other vectors
- Sort distances
- Return top K results



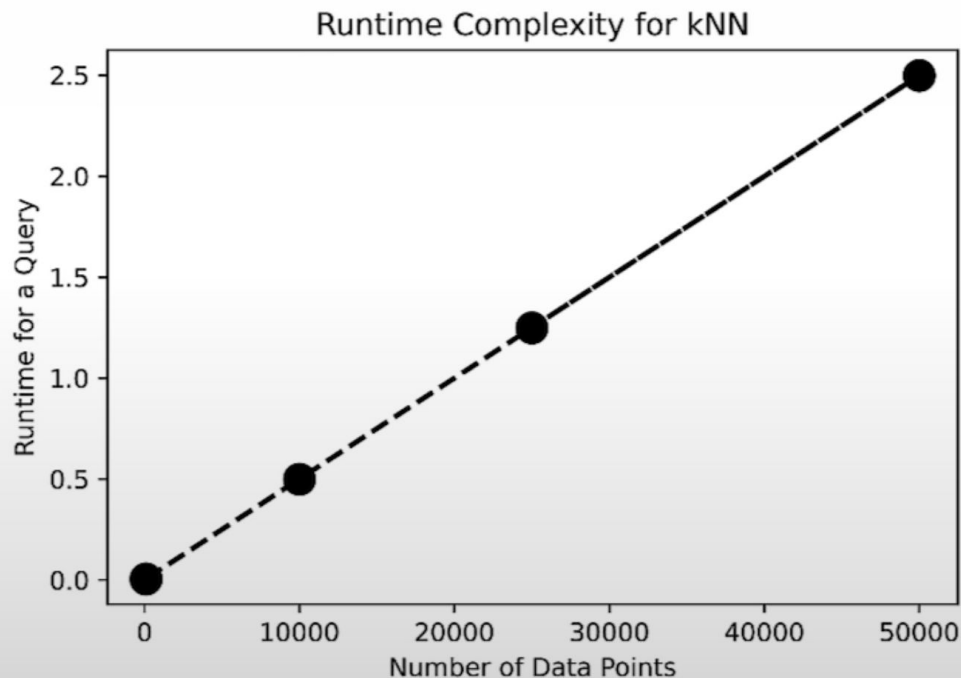




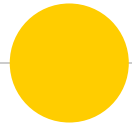




Performance



- Number of dimension
~1000
- Number of vectors
~millions
- Typical query latency
~seconds



Trade off accuracy for performance



Approximate Nearest Neighbour

- Preprocess data into efficient index
- Speed up search using index



Approximate Nearest Neighbour

- Tree-based algorithm (ANNOY)
- Locality sensitive hashing
- Product Quantization
- Inverted File Index
- **Proximity Graph (Hierarchical Navigable Small World)**



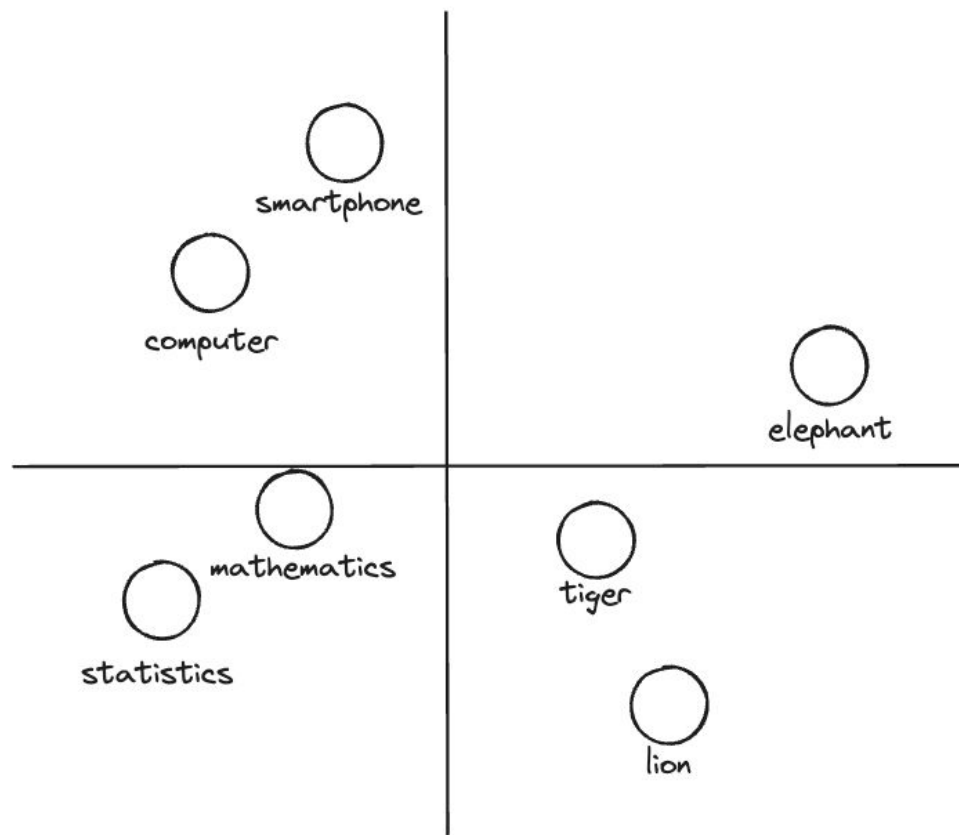
Navigable Small World

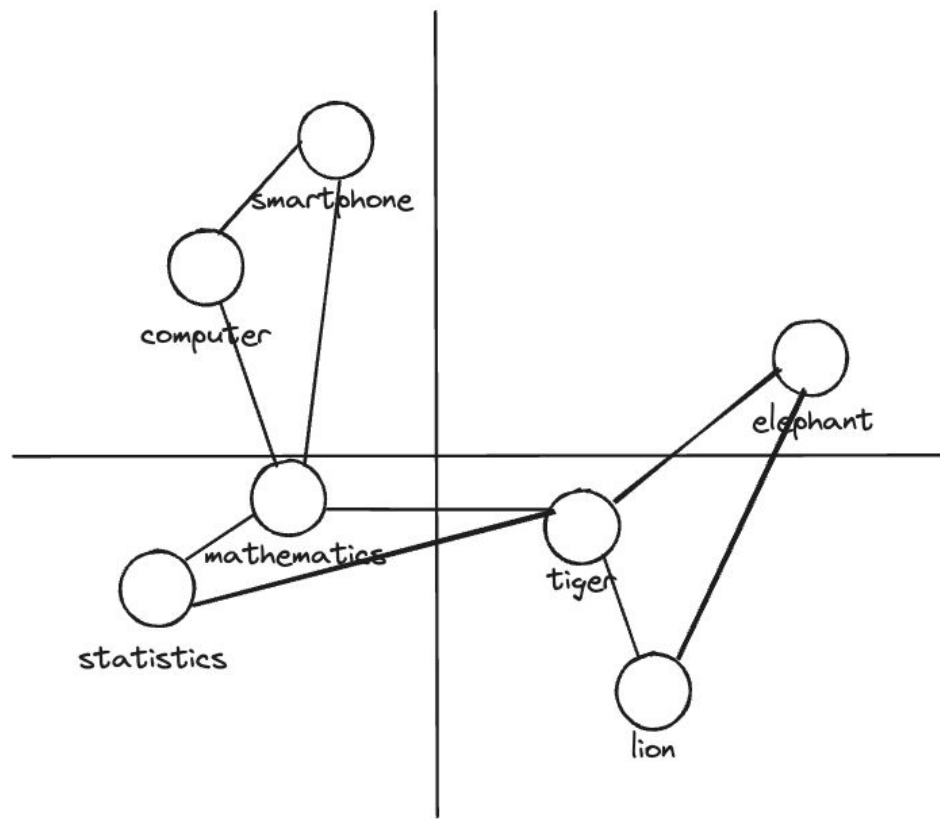
- Six degrees of separation
- Hungarian novelist – Karinthy Frigyes – 1929
- Better communication tech and travel
- Migrations and increasing connectivity
- "Shrinking" modern world

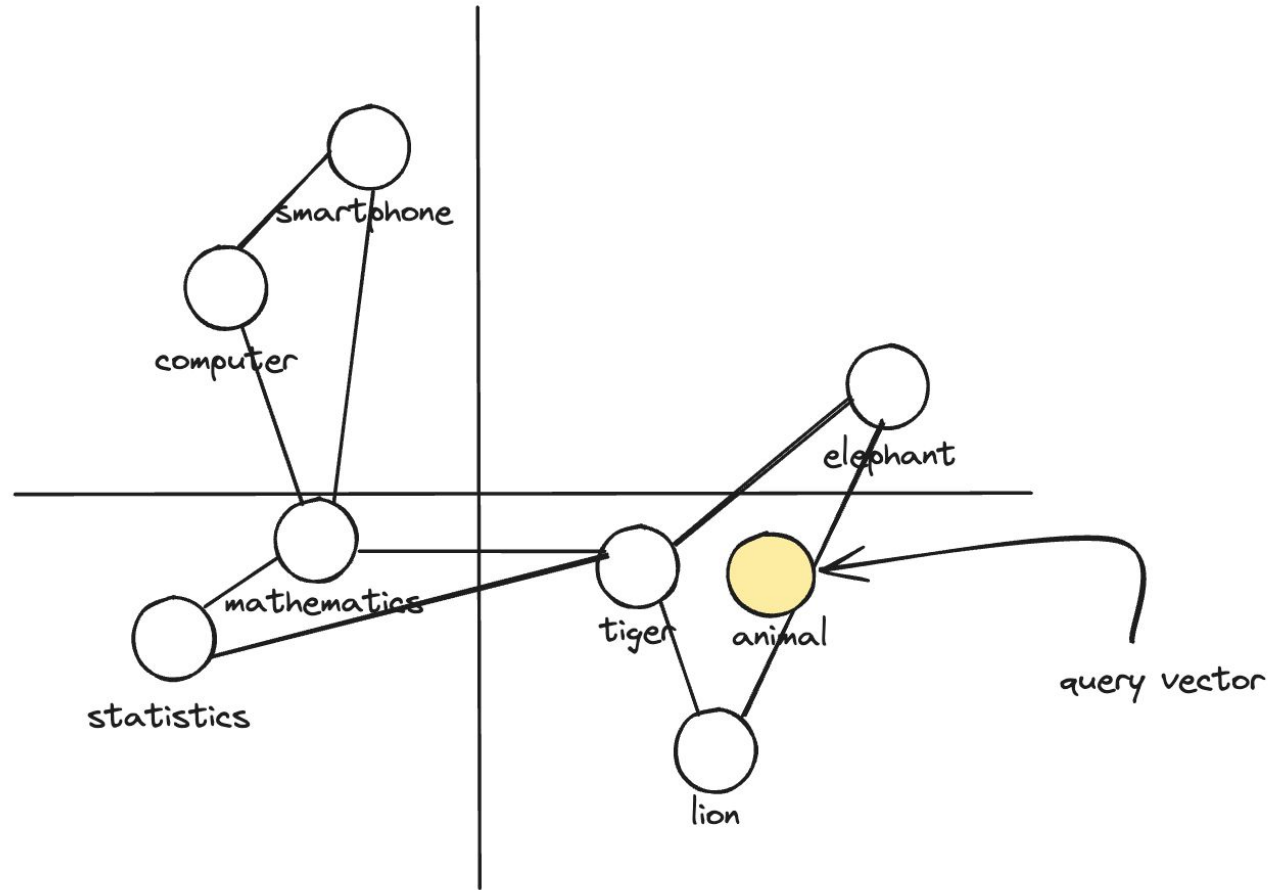


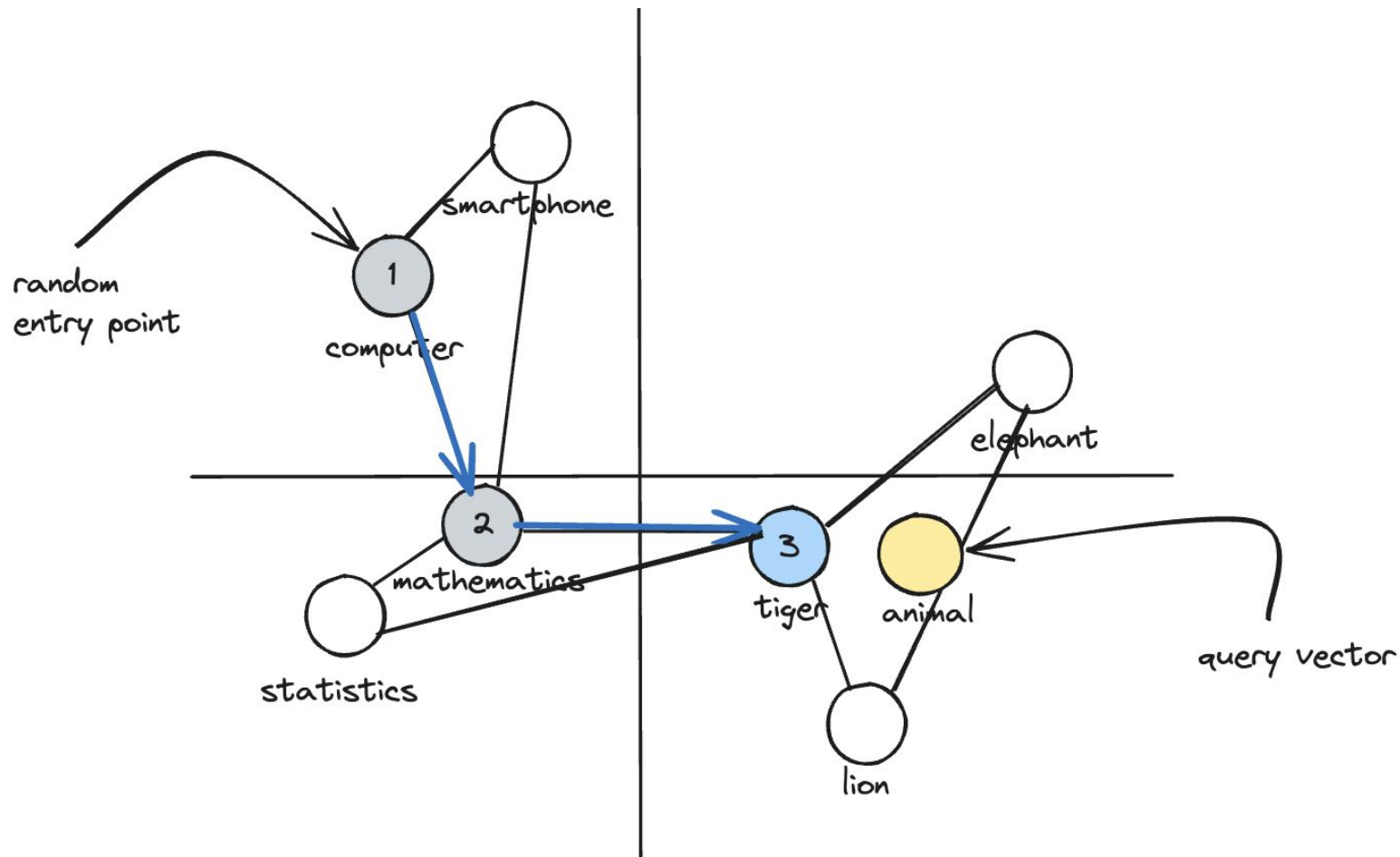
Navigable Small World

- every node can be reached with small number of hops from any other node
- Build a graph
- Every node has M connection to closest nodes (proximity list, “friend list”)
- Random entry node
- Greedy search, move to closest node
- Stop condition, no closer node in friend list





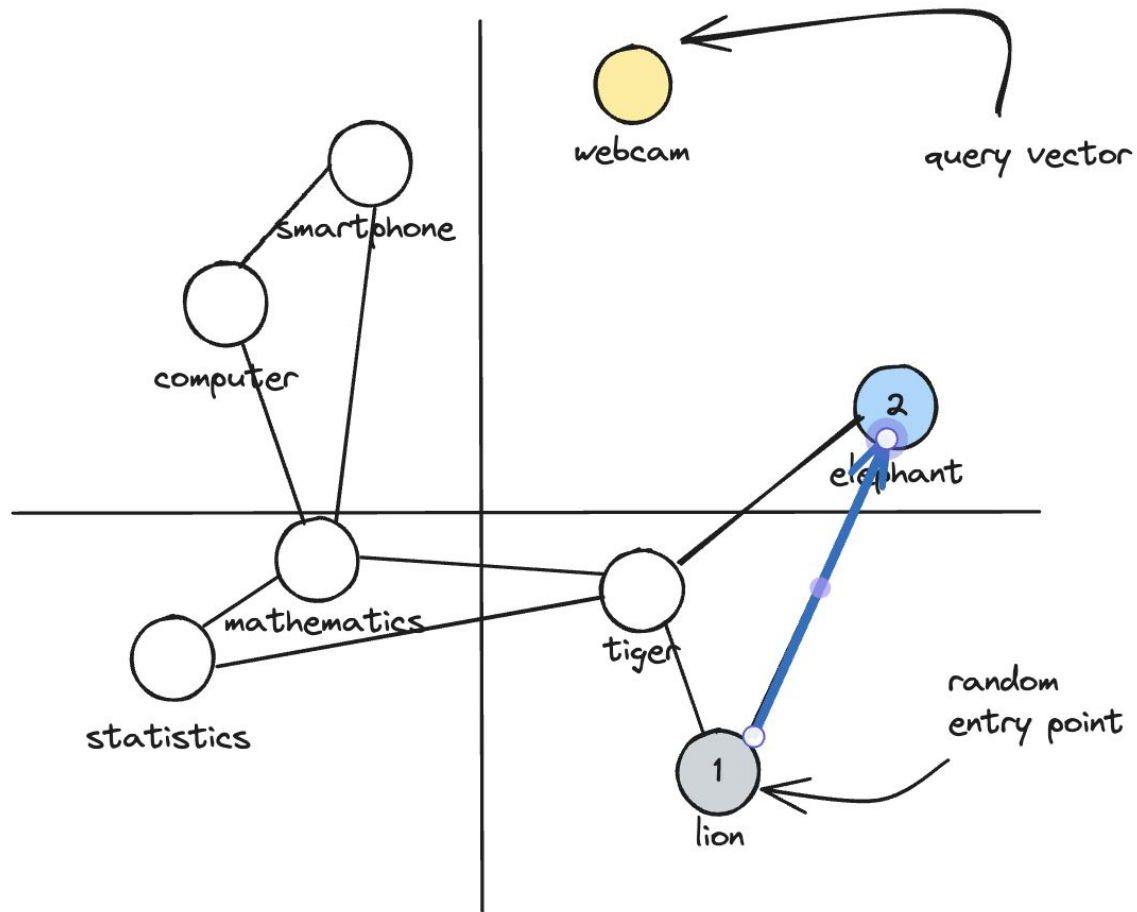






Drawbacks

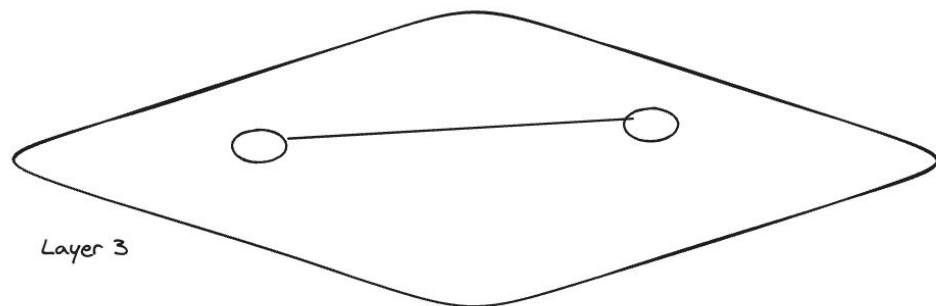
- Can get stuck in local minimum
- Polylogarithmic complexity, not performant at scale
 - Searches for multiple nodes, multiple times



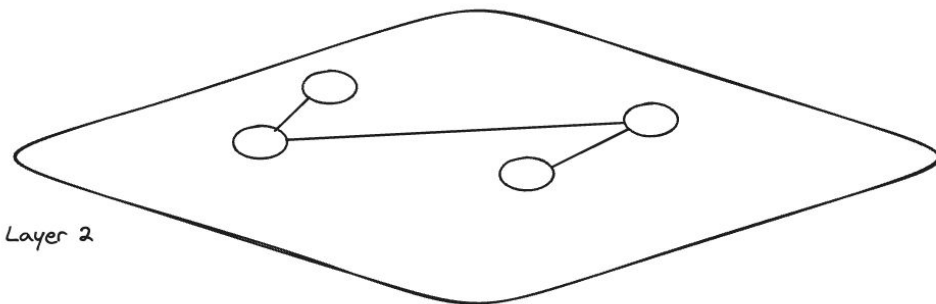


Hierarchical NSW

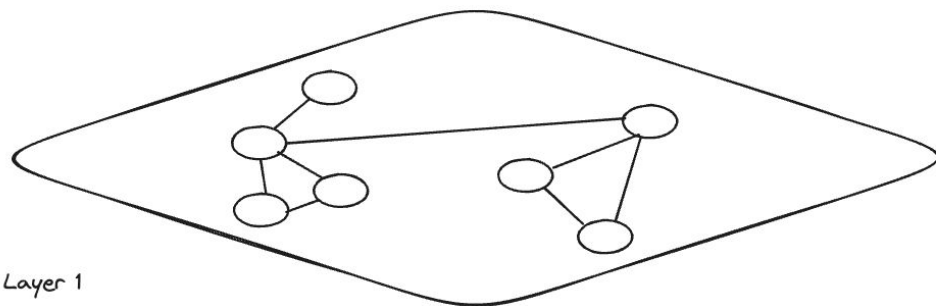
- Multi-layered graph
 - Base layer has all nodes (dense), more connections
 - Higher layer have fewer nodes (sparse) and fewer connections
- Search starts in top layer descending towards bottom layers



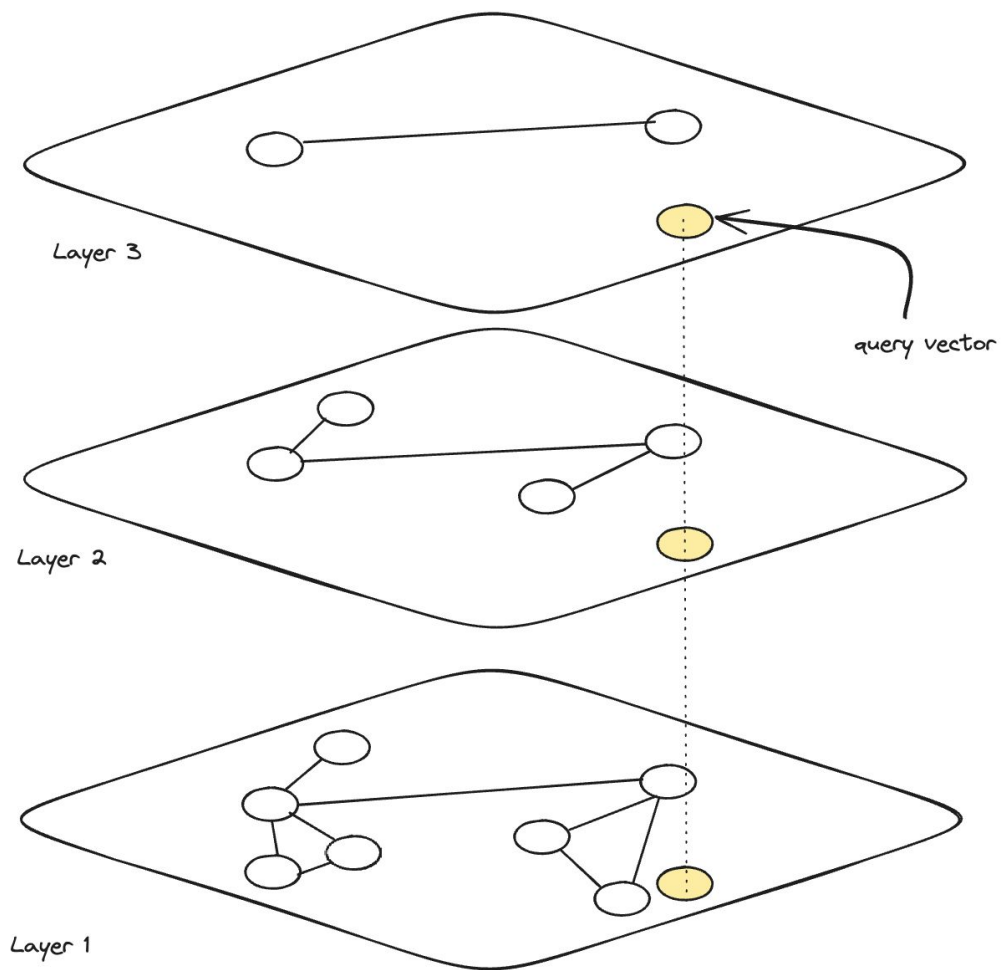
Layer 3

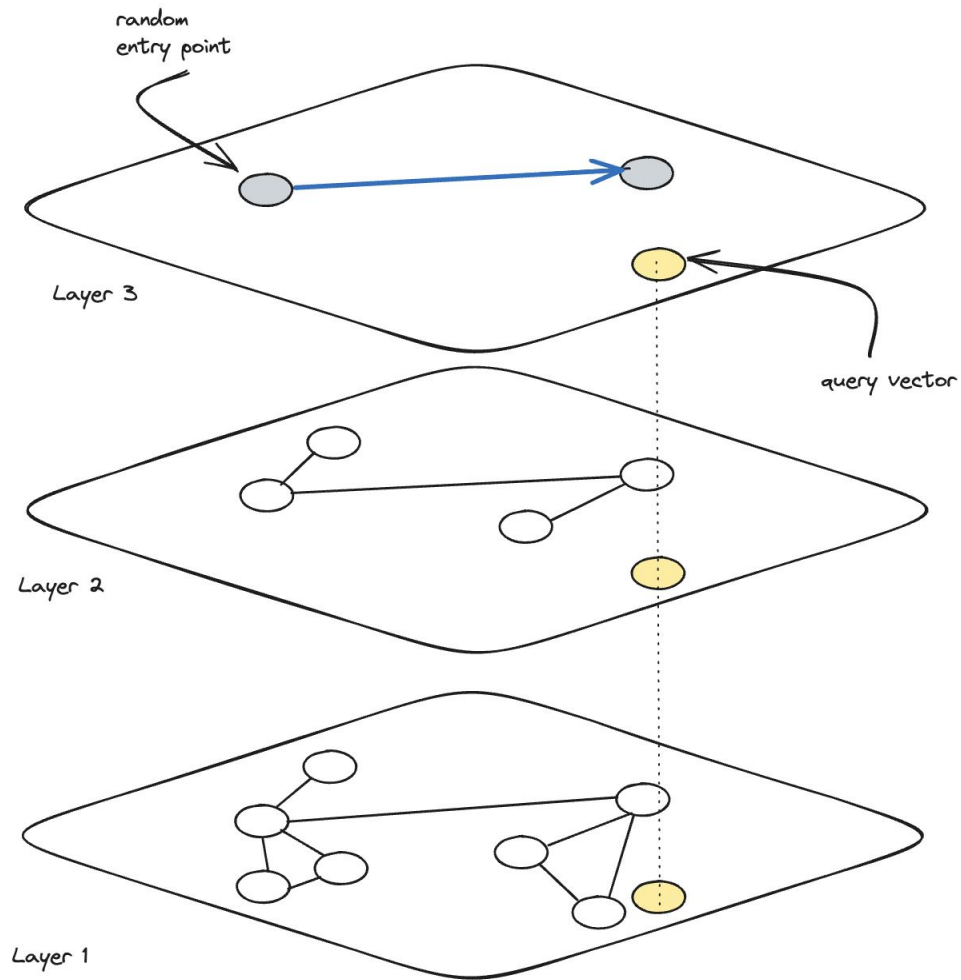


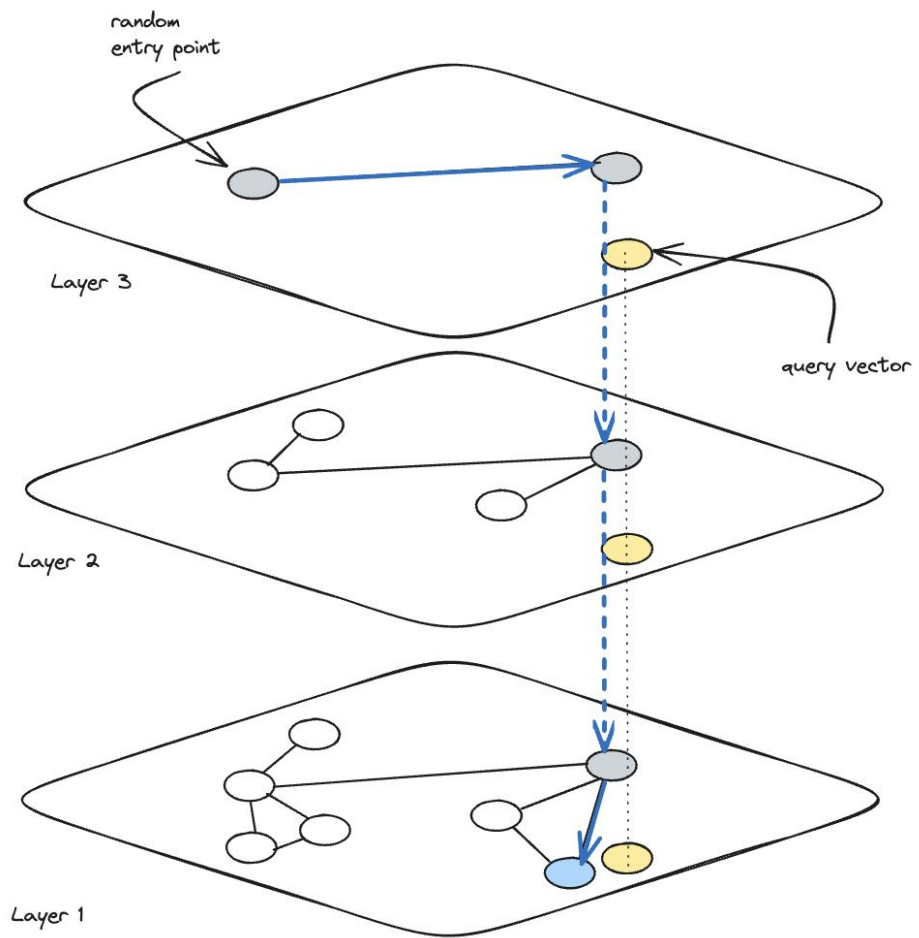
Layer 2



Layer 1

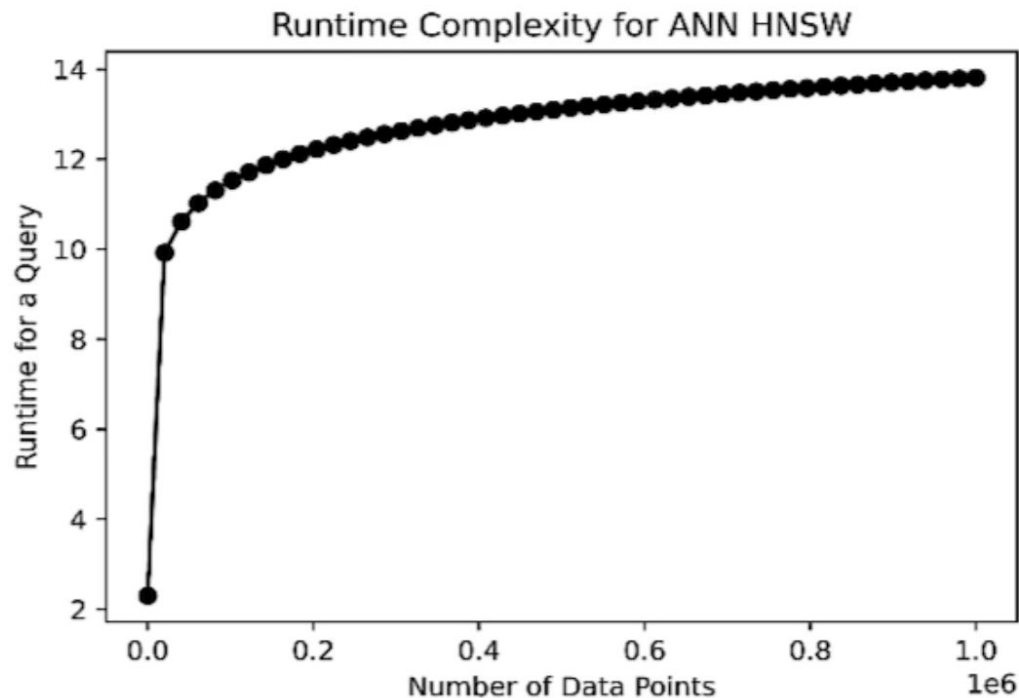








Performance



- Ensures consistent, reliable performance
- Query time increases logarithmically
- Typical query latency ~milliseconds