

EMBEDDINGS MASTERY SHEET

For Stephen Conversation - Monday Evening

CORE CONCEPT (Say This First)

"An embedding is a high-dimensional vector representation of text where semantic similarity becomes mathematical distance. Words or sentences with similar meanings end up close together in vector space, which lets us use math operations like cosine similarity to measure relatedness."

THE TWO KEY PROPERTIES

1. Semantic Representation

- Similar concepts → nearby vectors
- "cat" and "dog" are closer than "dog" and "strawberry"
- **Why it matters:** Captures meaning, not just keywords

2. Dimensionality

- **Lower dimensions (e.g., 128):** Faster, less memory, but less nuanced
- **Higher dimensions (e.g., 768, 1536, 7168):** Captures intricate relationships, but slower and can overfit
- **GPT-2:** 768 dimensions minimum
- **Modern models:** 1536+ dimensions

What "768 dimensions" means:

- Each word/sentence becomes a point with 768 coordinates
 - Like a map, but instead of (x, y), you have ($x_1, x_2, x_3, \dots, x_{768}$)
 - More coordinates = more precise location = more nuanced meaning
-

HOW SENTENCES BECOME VECTORS

You asked: "I thought only words can be vectorized?"

Method 1: Averaging (Simple)

"Help children learn"

 └ "help" → [0.2, 0.5, -0.3, ...]

 └ "children" → [0.4, 0.6, -0.1, ...]

 └ "learn" → [0.3, 0.7, -0.2, ...]

Average → [0.3, 0.6, -0.2, ...] = sentence embedding

Method 2: Dedicated Sentence Models (Better)

- Models like `sentence-transformers` or `text-embedding-3-small`
- Trained specifically to understand full sentences
- Captures context, not just word averages
- Example: "bank" in "river bank" vs "Chase bank" → different vectors

COSINE SIMILARITY EXPLAINED

Non-Technical: "How much do these vectors point in the same direction?"

The Math (Visual):

Vector A: → (pointing northeast)

Vector B: → (pointing northeast)

Similarity: 0.95 (very similar)

Vector A: → (pointing northeast)

Vector C: ← (pointing southwest)

Similarity: -0.95 (opposite)

Vector A: → (pointing northeast)

Vector D: ↑ (pointing north)

Similarity: 0.70 (somewhat similar)

Result range:

- 1.0 = identical direction (perfect match)
- 0.0 = perpendicular (unrelated)
- -1.0 = opposite direction (antonyms)

Why cosine and not just distance?

- Cosine ignores magnitude, focuses on direction

- "I love cats" and "I absolutely adore felines!!!" have different lengths but same meaning
 - Cosine captures this; Euclidean distance doesn't
-

YOUR NONPROFIT USE CASE

What You Built:

```
python

# 200 mission statements to compare
for mission_a, mission_b in all_pairs:
    prompt = f"Compare these missions and rate similarity 1-10:\n{mission_a}\n{mission_b}"
    response = llm_api_call(prompt) # Expensive, slow
```

Cost: 200 calls \times \$0.01 = \$2.00

Time: ~2 minutes

Problem: Using LLM reasoning for a similarity calculation

What You Should Have Built:

```
python

# ONE-TIME SETUP
embeddings = {}
for mission in all_missions:
    embeddings[mission.id] = embedding_api_call(mission.text) # Cheap, fast
    # Store in vector database

# INSTANT COMPARISONS
def find_similar(mission_id, top_n=10):
    query_vector = embeddings[mission_id]
    results = vector_db.similarity_search(query_vector, limit=top_n)
    return results # Returns in milliseconds

# OPTIONAL: Only use LLM for explanation
def explain_similarity(mission_a, mission_b, similarity_score):
    prompt = f"These missions have {similarity_score} similarity. Explain why in 2 sentences."
    return llm_api_call(prompt) # Only when needed
```

Cost: 200 embeddings \times \$0.0001 = \$0.02

Time: ~5 seconds

Benefit: 100x faster, 1/100th the cost, enables new features

NEW FEATURES EMBEDDINGS ENABLE

1. Instant "Find Similar" Search

- User searches: "education nonprofits"
- Embed the query → find nearest vectors → return results
- No keyword matching needed!

2. Clustering

- Automatically group similar missions
- "These 50 nonprofits all focus on youth education"

3. Anomaly Detection

- "This mission statement is very different from others in this category"

4. Smart Recommendations

- "Users who supported Mission A also supported these similar missions..."
-

TOKEN EMBEDDINGS vs CONTEXTUAL REPRESENTATIONS

Important distinction for technical depth:

Token Embeddings (Static)

- Fixed vectors assigned at input layer
- "bank" always gets the same initial vector
- **These are what we're talking about for your use case**

Contextual Representations (Dynamic)

- Produced by deeper layers of the model
- "bank" gets different representation in "river bank" vs "Chase bank"
- Change as they flow through transformer layers
- Also called "hidden states" or "contextual embeddings"

For Stephen: "When I talk about embeddings for the nonprofit feature, I mean static sentence embeddings from a dedicated embedding model, not the contextual representations from LLM hidden layers."

EVOLUTION OF EMBEDDINGS (Quick History)

Traditional (Statistical)

- **TF-IDF:** Based on word frequency
- **Problem:** "cat" and "dog" have no relationship
- **Use:** Keyword search, document classification

Static Word Embeddings

- **Word2Vec, GloVe:** Capture semantic relationships
- **Problem:** "bank" always has same vector regardless of context
- **Use:** Recommendation systems, similarity search

Contextual Embeddings (Modern)

- **BERT, GPT family:** Context-aware representations
- **Benefit:** "bank" adapts to sentence meaning
- **Use:** Everything modern LLMs do

For your nonprofit case: Static sentence embeddings are perfect (fast, cheap, effective)

COMMON INTERVIEW QUESTIONS

Q: "Why use embeddings instead of keyword matching?"

A: "Keyword matching only finds exact words. Embeddings capture semantic meaning, so 'youth education' matches 'teaching children' even with no shared words."

Q: "How do you choose embedding dimensions?"

A: "It's a tradeoff. I'd start with a pre-trained model like OpenAI's text-embedding-3-small (1536 dims) - it's optimized for quality and speed. Only fine-tune or adjust if I had specific performance needs."

Q: "What's the difference between embedding and fine-tuning?"

A: "Embeddings convert text to vectors for similarity/search. Fine-tuning retrains model weights for better task-specific performance. For nonprofit comparison, embeddings are the right tool - we're solving a similarity problem, not teaching the model new behavior."

WHEN STEPHEN SAYS "YOU SHOULD KNOW EMBEDDINGS"

Your response: "You're absolutely right, and I've been digging into this since our conversation. I realized we were using LLM reasoning for what's fundamentally a similarity problem. The better architecture would be embedding all missions once, storing them in a vector database, and using cosine similarity for comparisons."

We'd only call the LLM when we need qualitative explanation. It would be 100x faster and cost 1/100th as much. That's a lesson I've taken to heart - matching the right tool to the problem type."

TECHNICAL TERMS TO USE CONFIDENTLY

- **Vector space:** The mathematical space where embeddings exist
 - **Cosine similarity:** Angle-based similarity metric (direction, not distance)
 - **Dimensionality:** Number of coordinates in the vector
 - **Semantic representation:** Capturing meaning mathematically
 - **Sentence transformers:** Models specialized for sentence embeddings
 - **Vector database:** Database optimized for similarity search (e.g., Pinecone, Weaviate)
-

RED FLAGS TO AVOID

- ✗ "Embeddings are just like TF-IDF" (No - embeddings capture semantics)
 - ✗ "We fine-tune embeddings for our use case" (No - you use pre-trained embedding models)
 - ✗ "Each word is 768 numbers" (Close, but say "768-dimensional vector" or "point in 768-dimensional space")
 - ✗ Confusing embeddings with the full LLM
-

YOU'RE READY. You understand the concept, the application, and the business value. That's what Stephen wants to see.