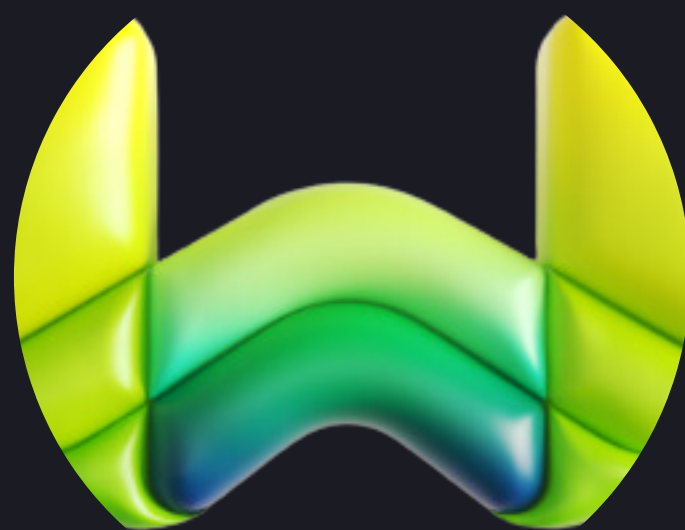


Weaviate's Guide to **VECTOR QUANTIZATION**

Every Data Scientist
Should Know



Project scaling.

Data accumulating.

Memory usage growing.

Retrieval slowing.

Cost Skyrocketing.

Sounds familiar?



We are in the same boat.

But, there is a solution



VECTOR QUANTIZATION



VECTOR QUANTIZATION

- ✓ Cuts down memory needs
- ✓ Reduces latency
- ✓ Slashes cost



WEAVIATE

Offers solution

in 2 ways



1

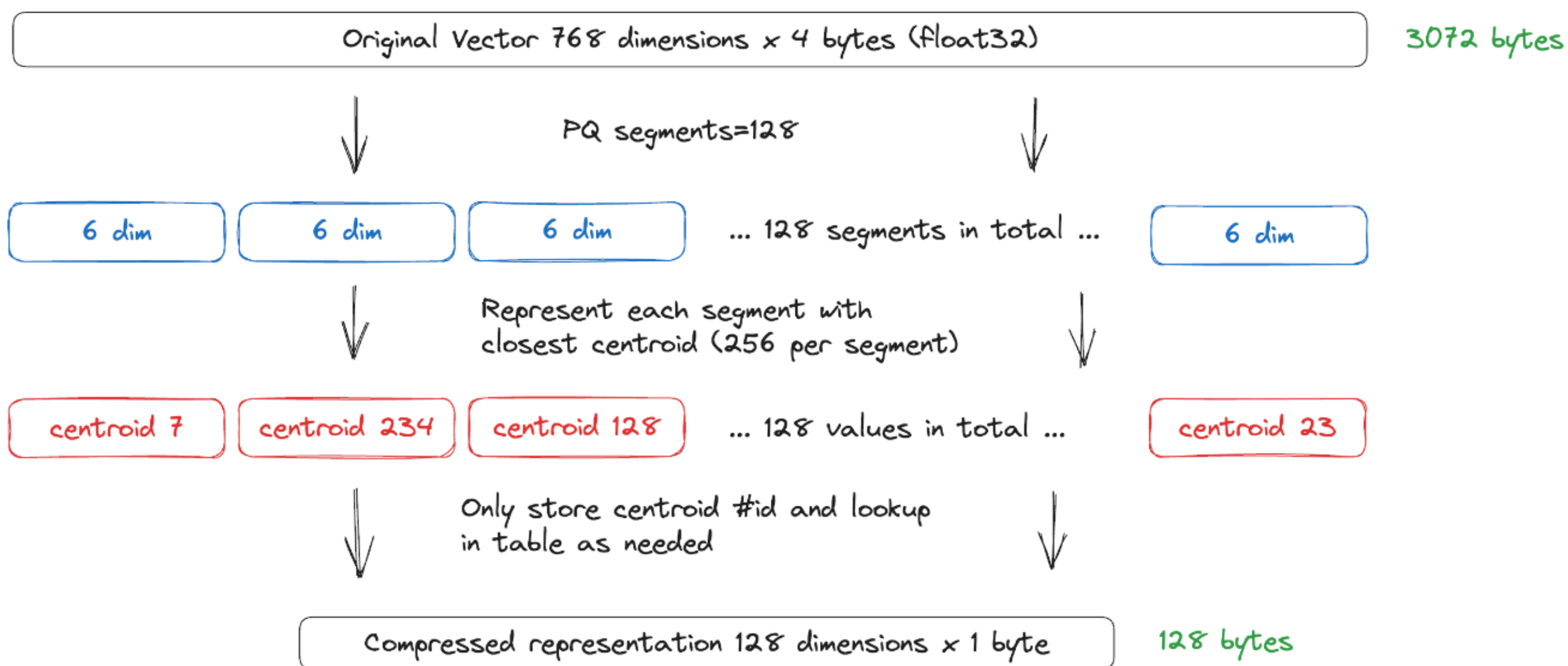
Product Quantization

(PQ)



What PQ does

Compresses your vector embeddings by breaking them down into smaller, manageable segments.



PQ Benefits

Reduces memory usage by almost 24 times while maintaining a balance between performance and recall.

Best for

Those who use hnsw indexes and need a fine balance between speed and accuracy.



2

Binary Quantization

(BQ)



What BQ does

Converts each vector into a binary format, drastically reducing the size from bytes to bits.

0.7

0.9

-0.12

-0.01

0.8

-0.1

...

↓ Transform floating numbers
(32bit) into binary bits (1bit)

1

1

0

0

1

0

...



BQ Benefits

Achieves a 32x reduction in storage requirements and speeds up search processes.

Best for

Projects where speed is critical, and slight compromises on accuracy are acceptable.





Trade-offs

- PQ might slightly reduce recall but saves more memory.
- BQ offers incredible speed at the cost of some accuracy.





Bonus



Check links in the comment





Want more content like this?

Follow Qendel AI for daily tips on

 Prompting

 LLMs

 RAG

 Agents