

Vector Length Analysis on Real Data + PCA + Arnon's Work

Distributed Monitoring

Table Of Contents

- ❑ Comparison to Arnon's results
- ❑ Bag of Words – Inner Product
- ❑ Bag of Words – Entropy
- ❑ PCA results

Comparison to Arnon's Work

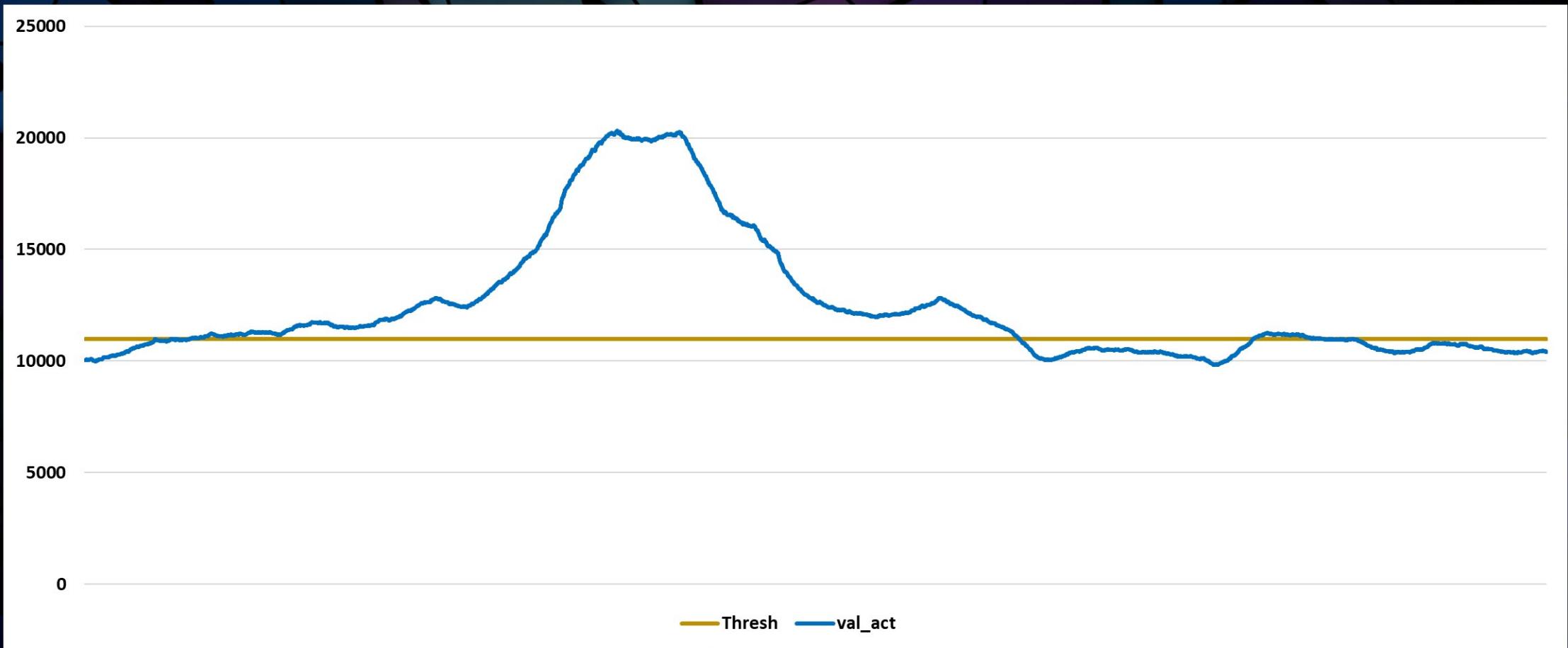
Tweeter Data Set

- 10 nodes – Round Robin
- Bag of words vector - $|\text{vector}| = 1254 \cdot 2 (X \cdot Y)$
- Constant threshold – 11,000
- Global BOW is the AVERAGE of BOWs
- Window size = 1,000
- Step – one tweet

Inner Product
Value

Arnon's Data

↑
Value
→ Time

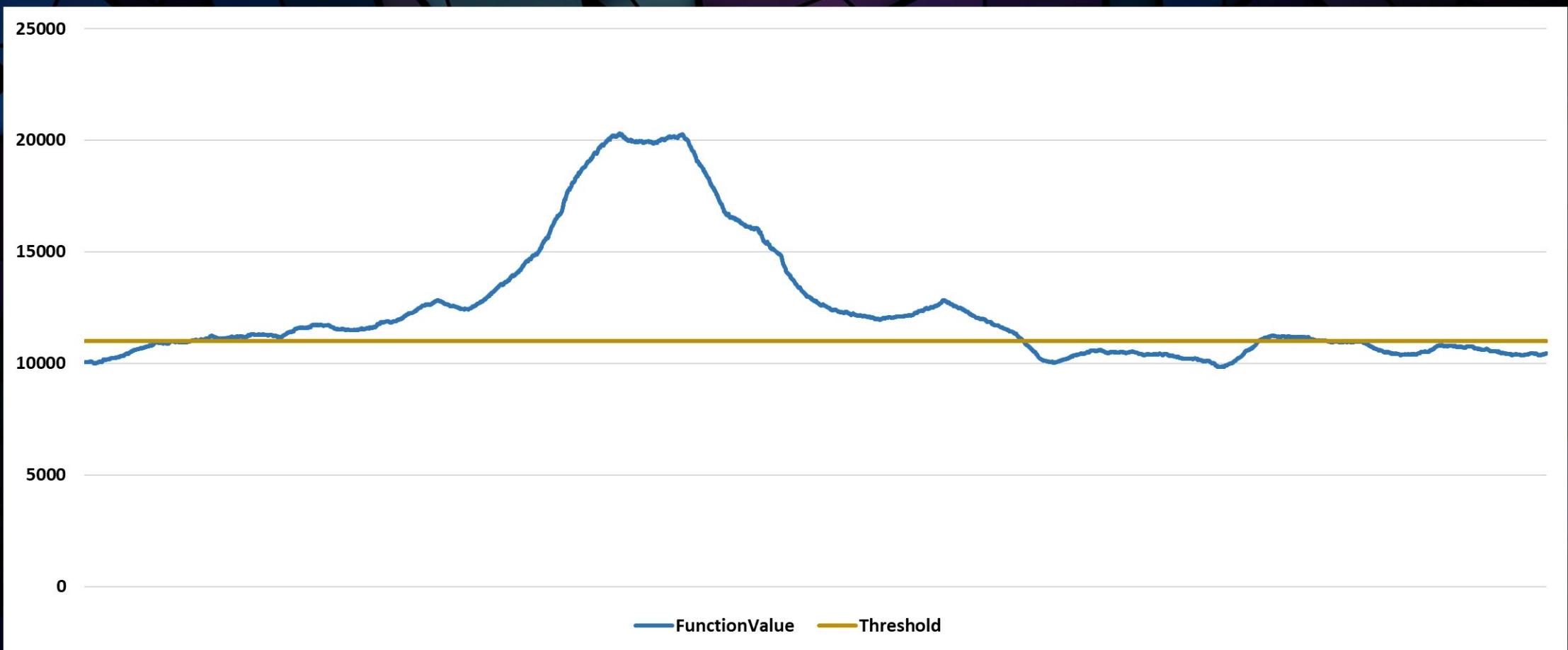


Inner Product
Value



Time

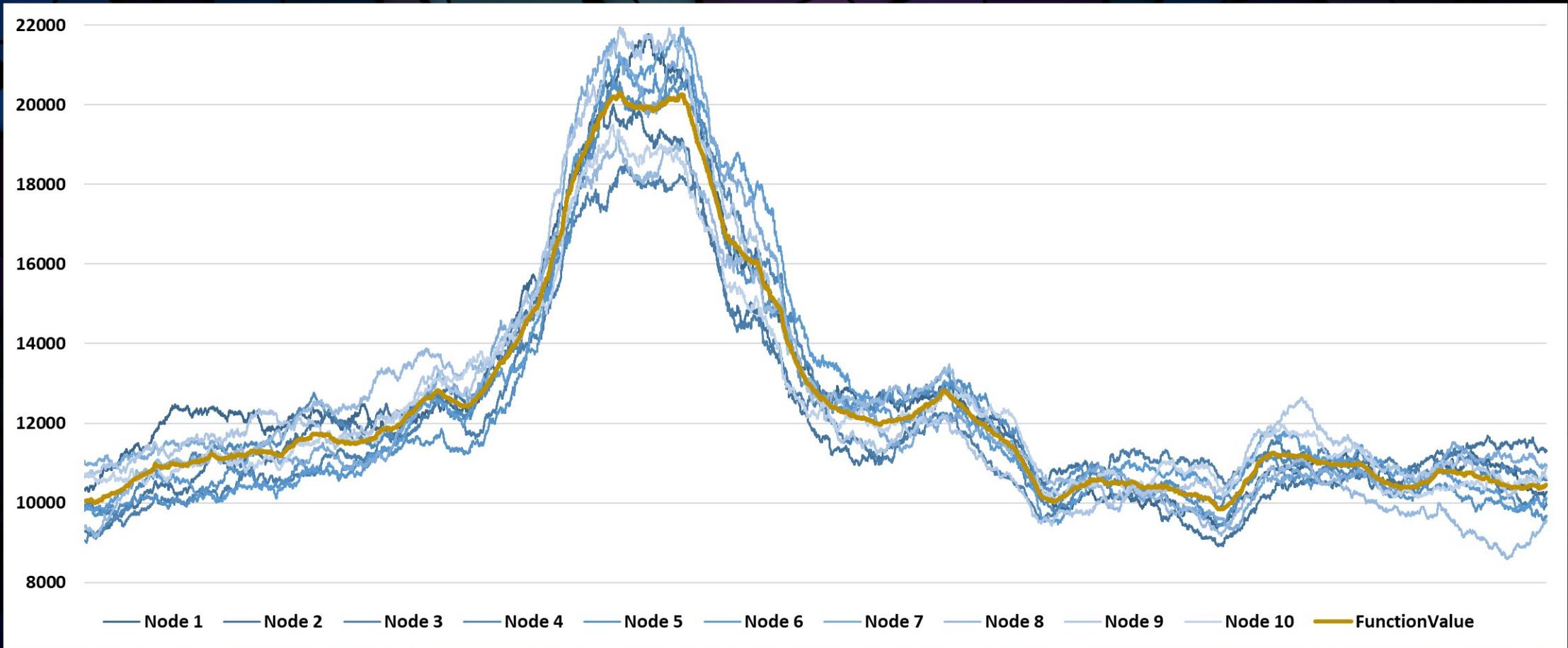
My Reconstructed Function



DATA

Inner Product
Value

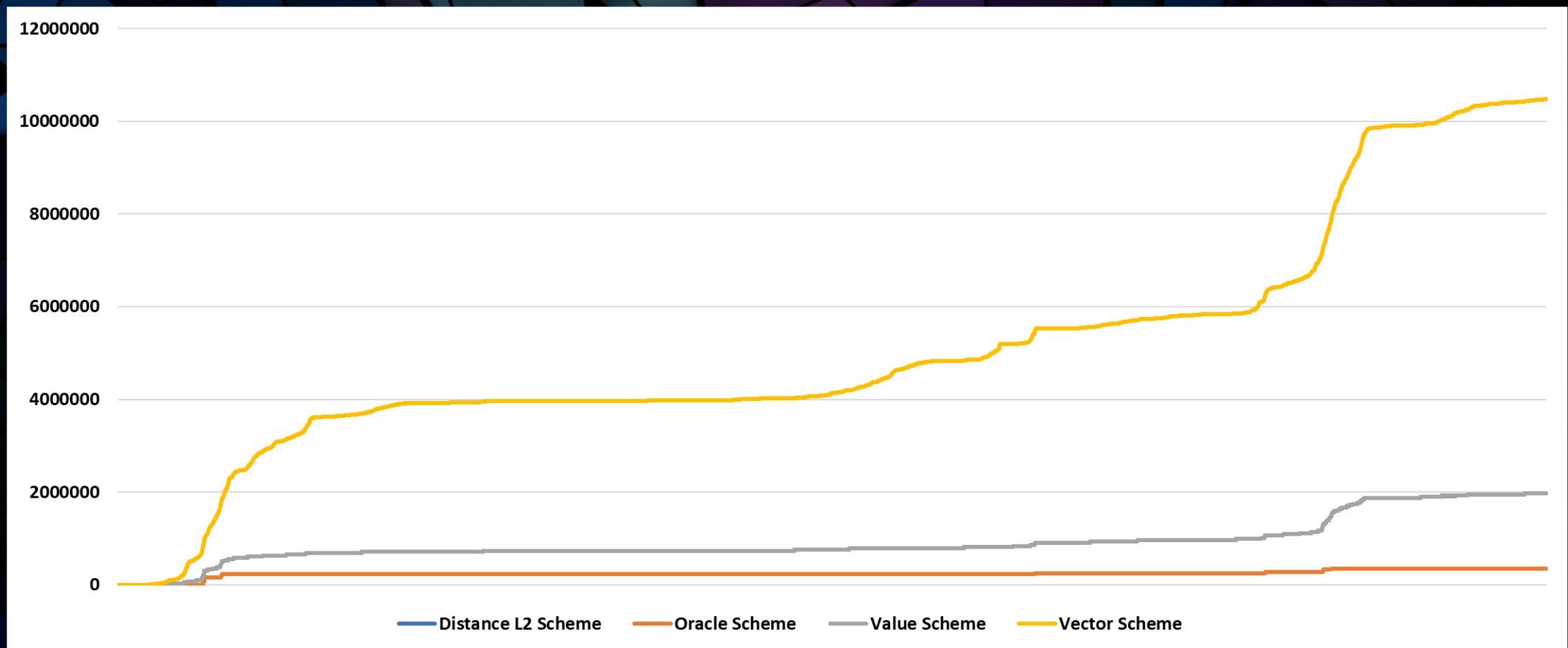
Time



Bandwidth

Time

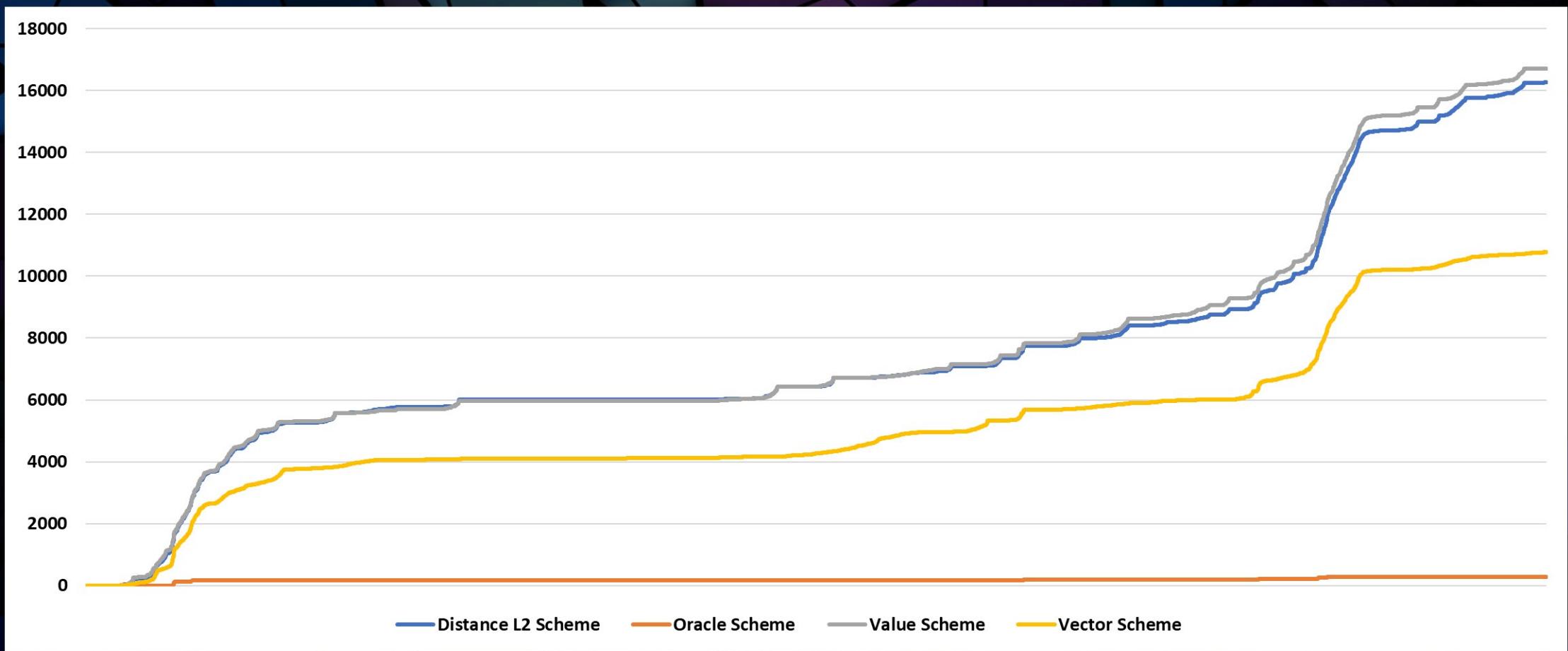
Bandwidth Over Time



Messages

Messages Over Time

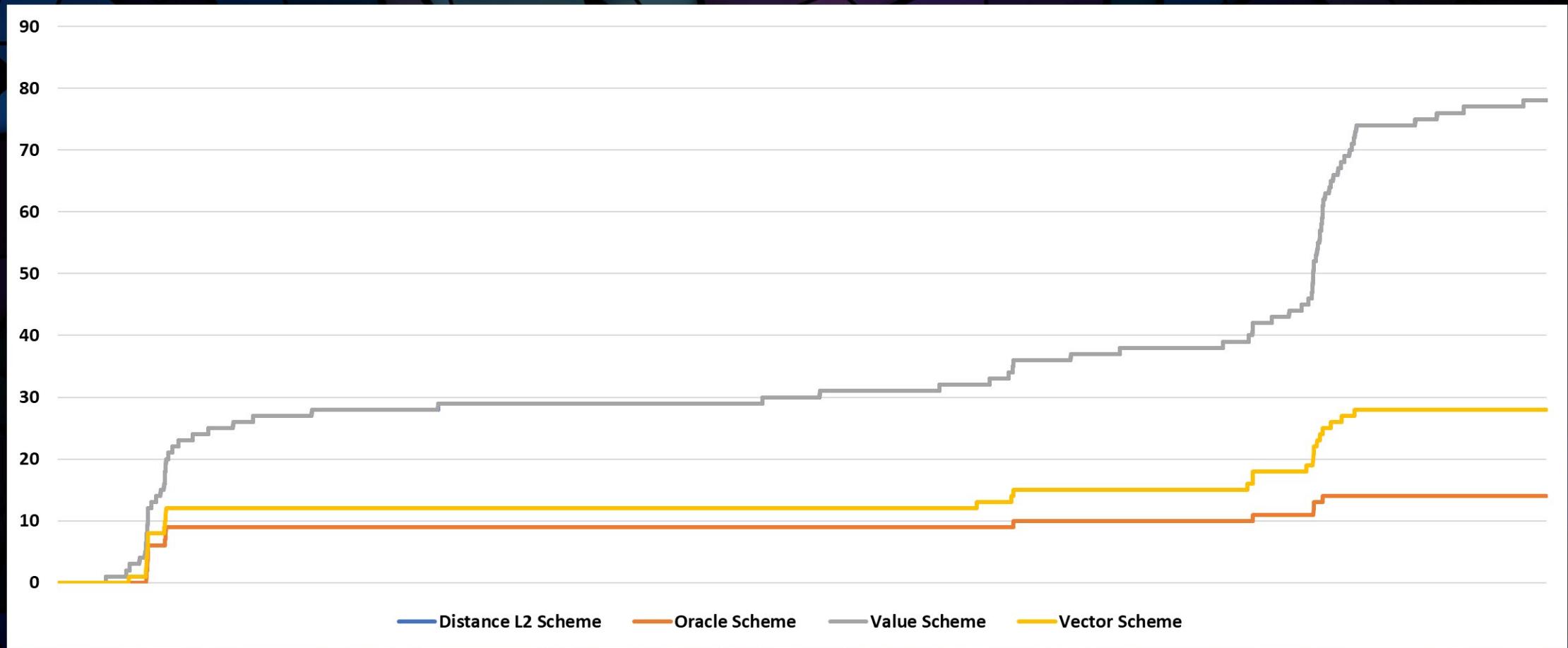
Time



Full Syncs

Full Syncs Over Time

Time

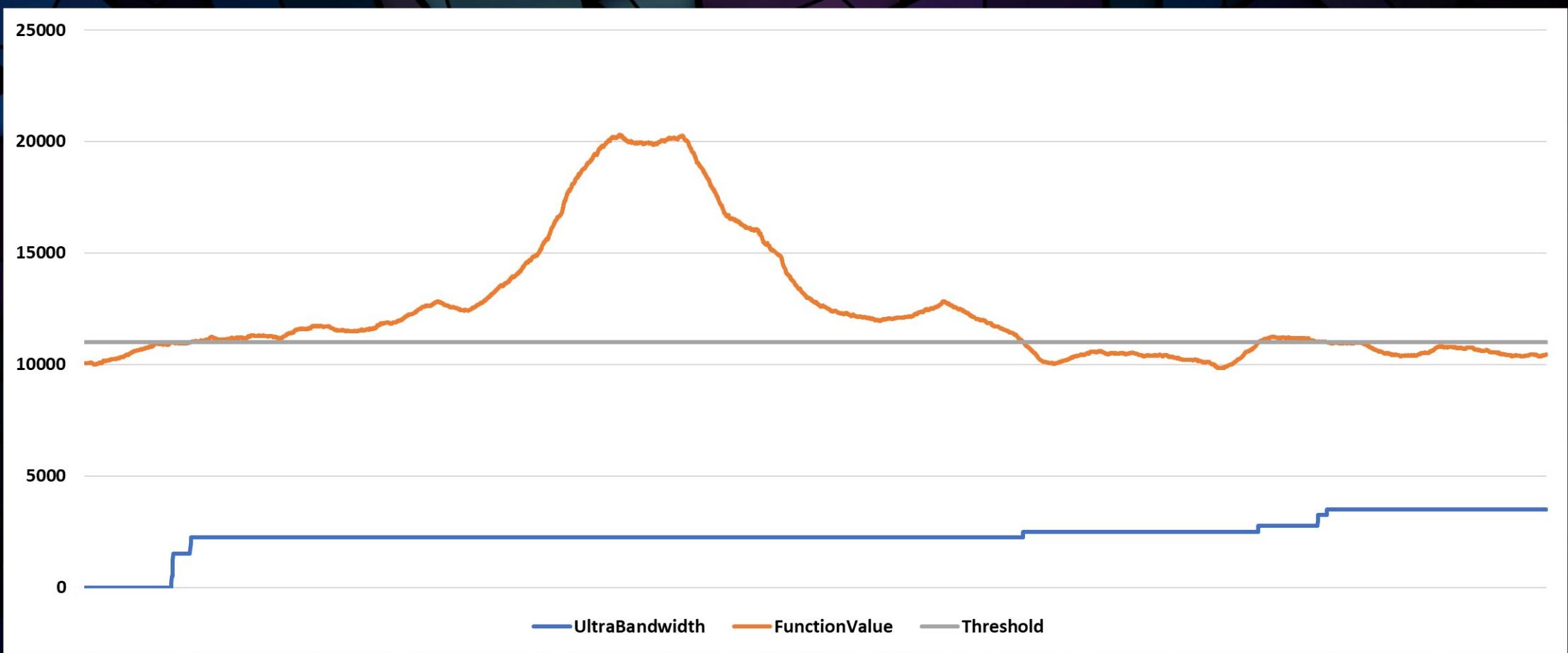


Bandwidth/
Value



Time

Oracle Bandwidth Scaled

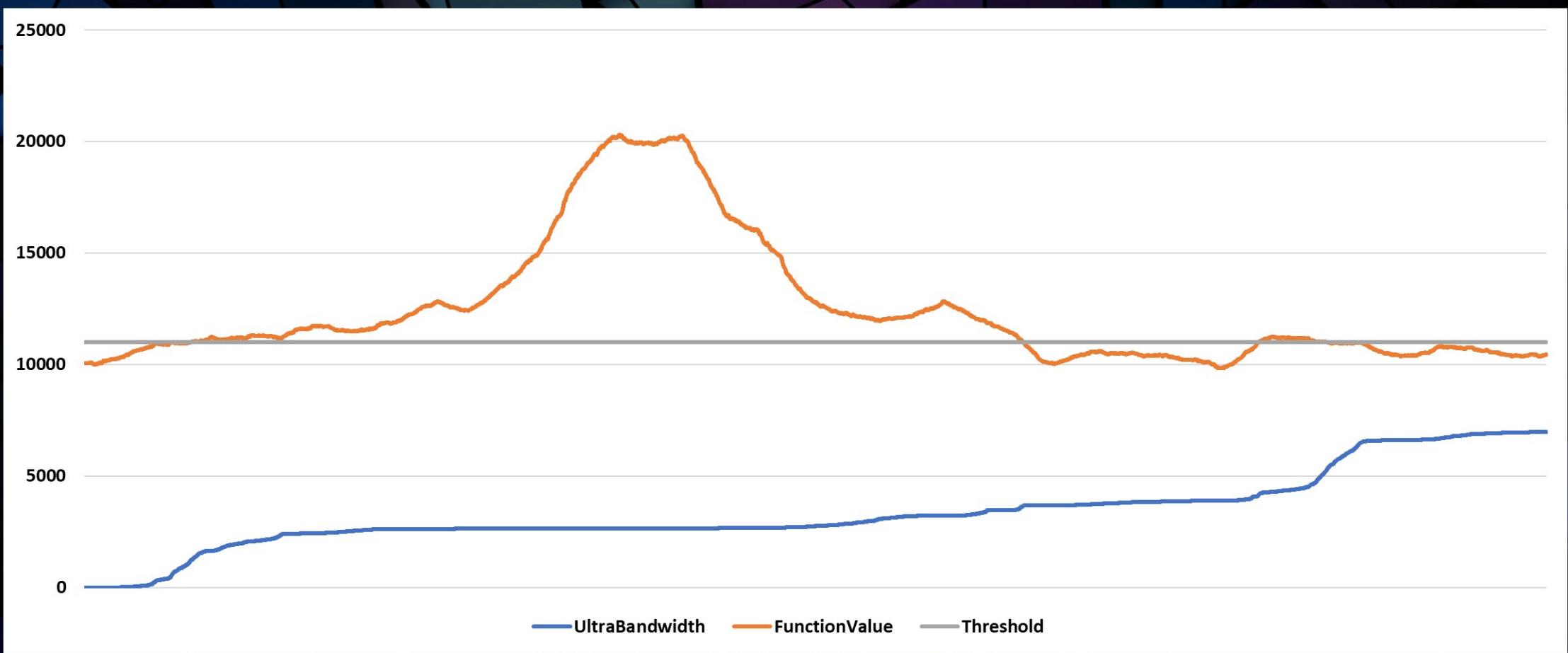


Bandwidth/
Value



Time

Vector Bandwidth Scaled

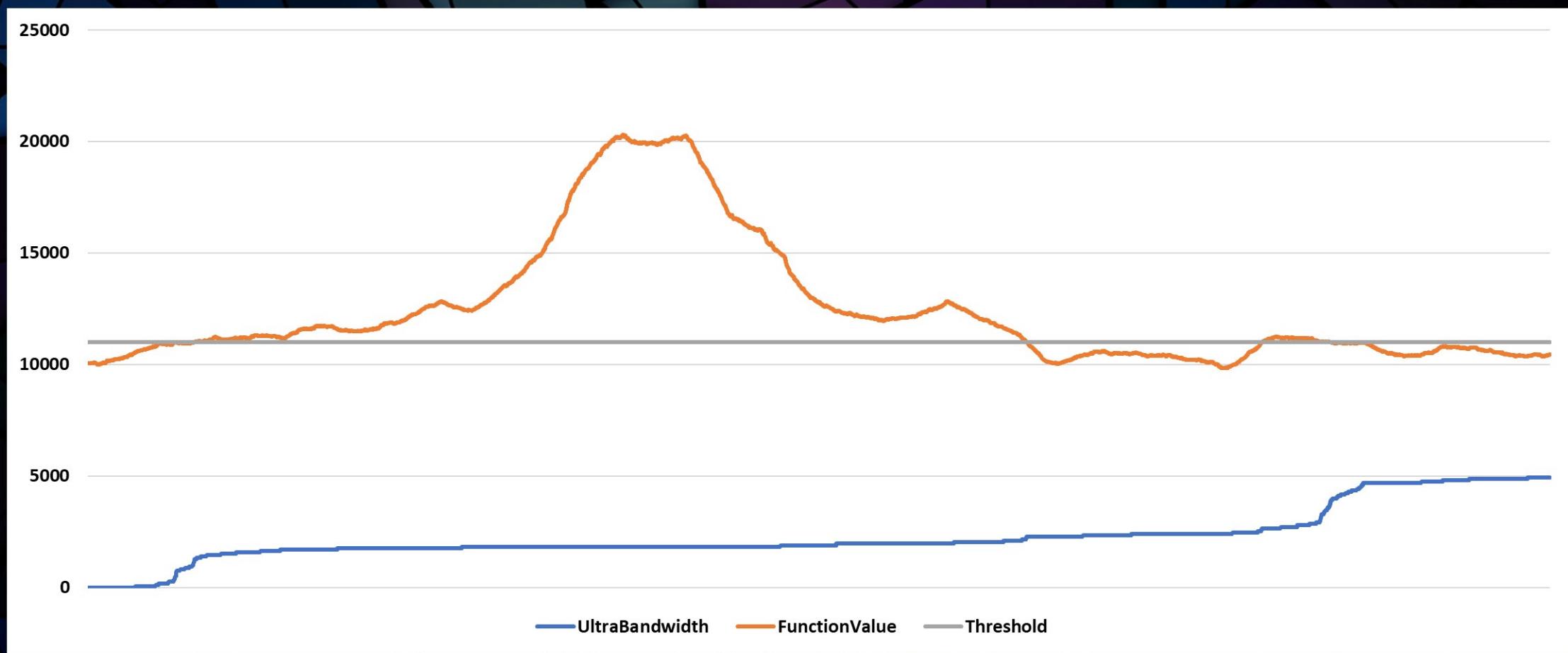


Bandwidth/
Value



Time

Distance Bandwidth Scaled



Bag of Words

Bag of Words

- ❑ Bag of words of length k consists of the **K MOST COMMON TOKENS** in English:
 - ❑ ('the', 'of', 'and', 'to', 'a', 'in', 'for', 'is', 'on'...)
 - ❑ tokens not in the k most common are ignored
 - ❑ a histogram of the occurrences moves in a '**sliding window**' fashion

Data Sources

- Textual data, each node is fed of a different oriented data:
 - Blogs textual data
(<http://u.cs.biu.ac.il/~koppel/BlogCorpus.htm>)
 - India's news headlines textual data
(<https://www.kaggle.com/therohk/india-headlines-news-dataset>)
 - Jeopardy questions textual data
(https://www.reddit.com/r/datasets/comments/1uyd0t/200000_jeopardy_questions_in_a_json_file/)
 - Spam clickbaits of news headlines textual data
(<https://www.kaggle.com/therohk/examine-the-examiner>)
 - Some books textual data

Inner Product

Inner Product Monitoring

- ❑ Distance Scheme – L_2 Distance to convex bound
- ❑ Data vector:

`vector(i) <- # occurrences in the window of the i'th lexicographic
ordered n-most common token`

`Inner Product(vector) = vector[0..n/2] · vector[n/2...n]`

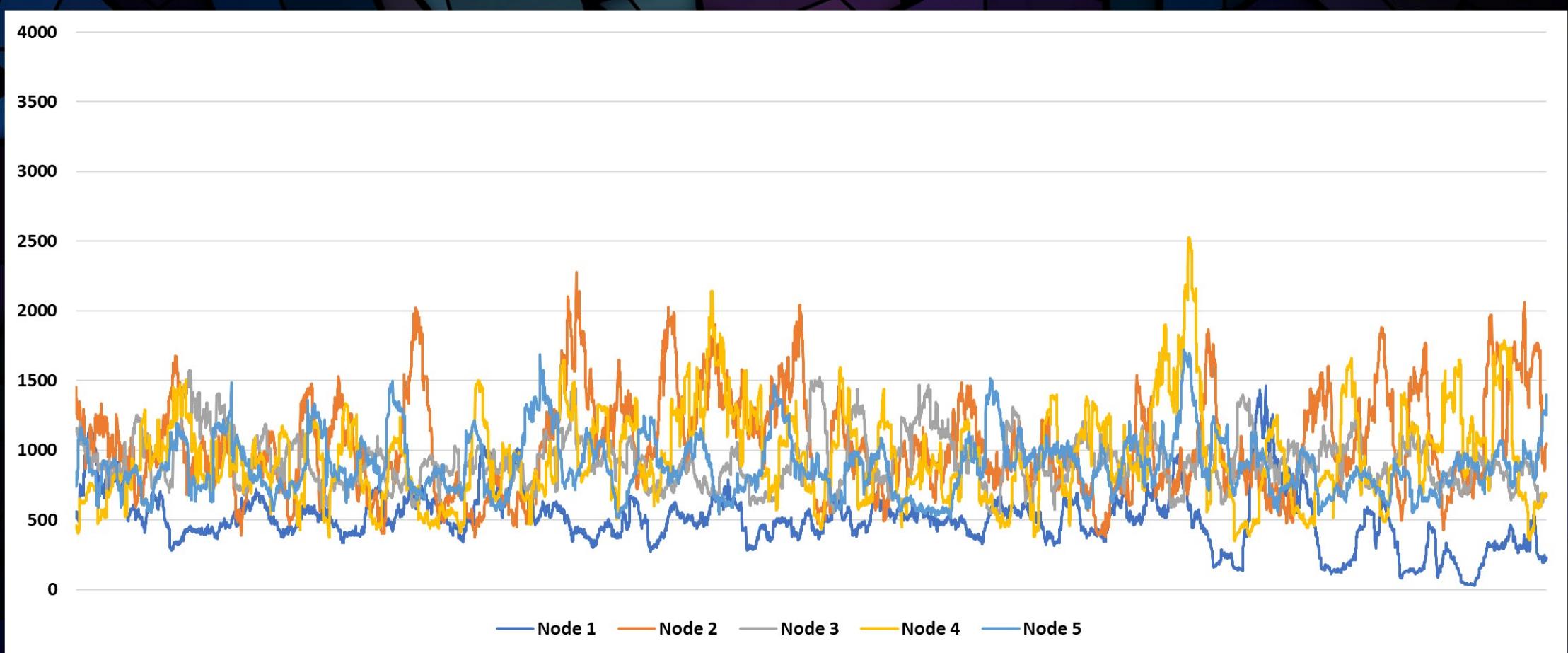
Running Information

- ❑ Window Size = 50,000
- ❑ Step Size = 100
- ❑ Epsilon – Multiplicative (*0.988, *1.012)
- ❑ Global BOW is the SUM of the local BOWs

Inner
Product

Nodes, Inner Data, $|\text{vector}| = 10000$

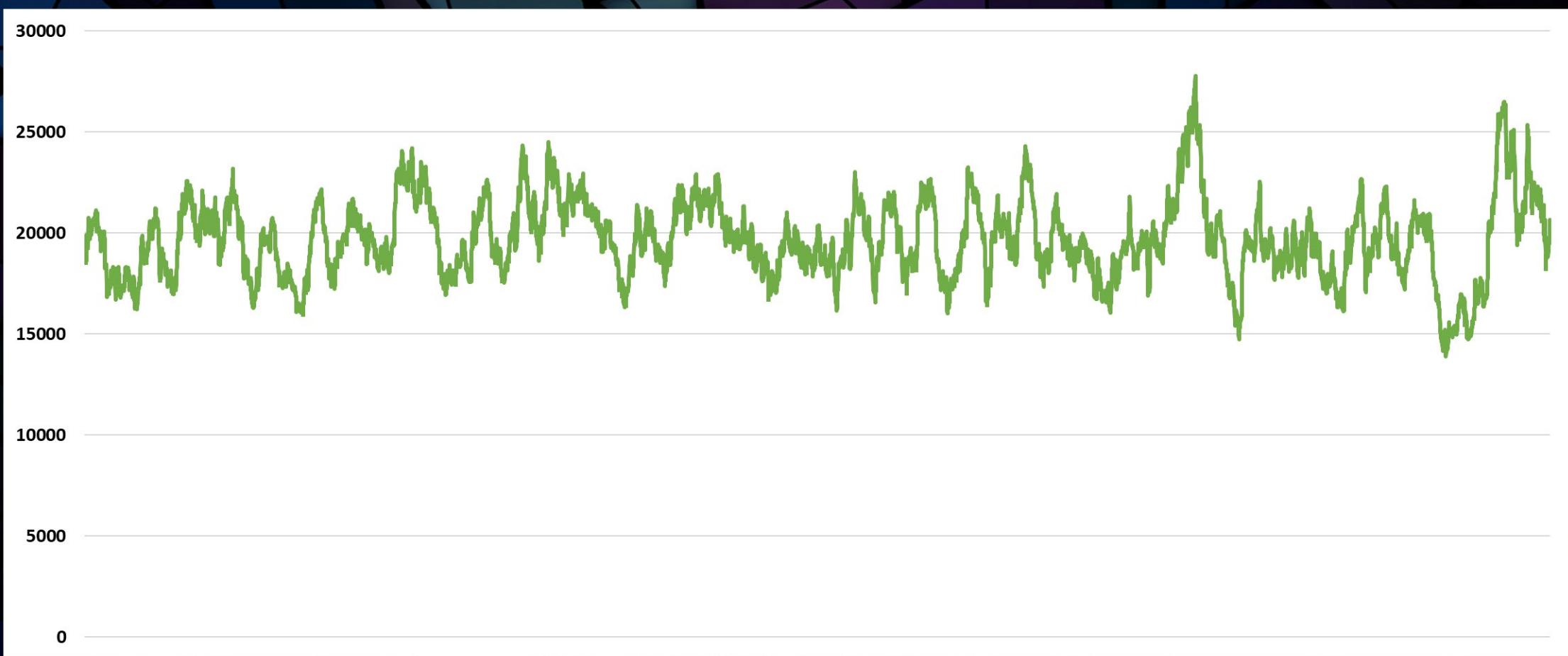
Time



Inner
Product

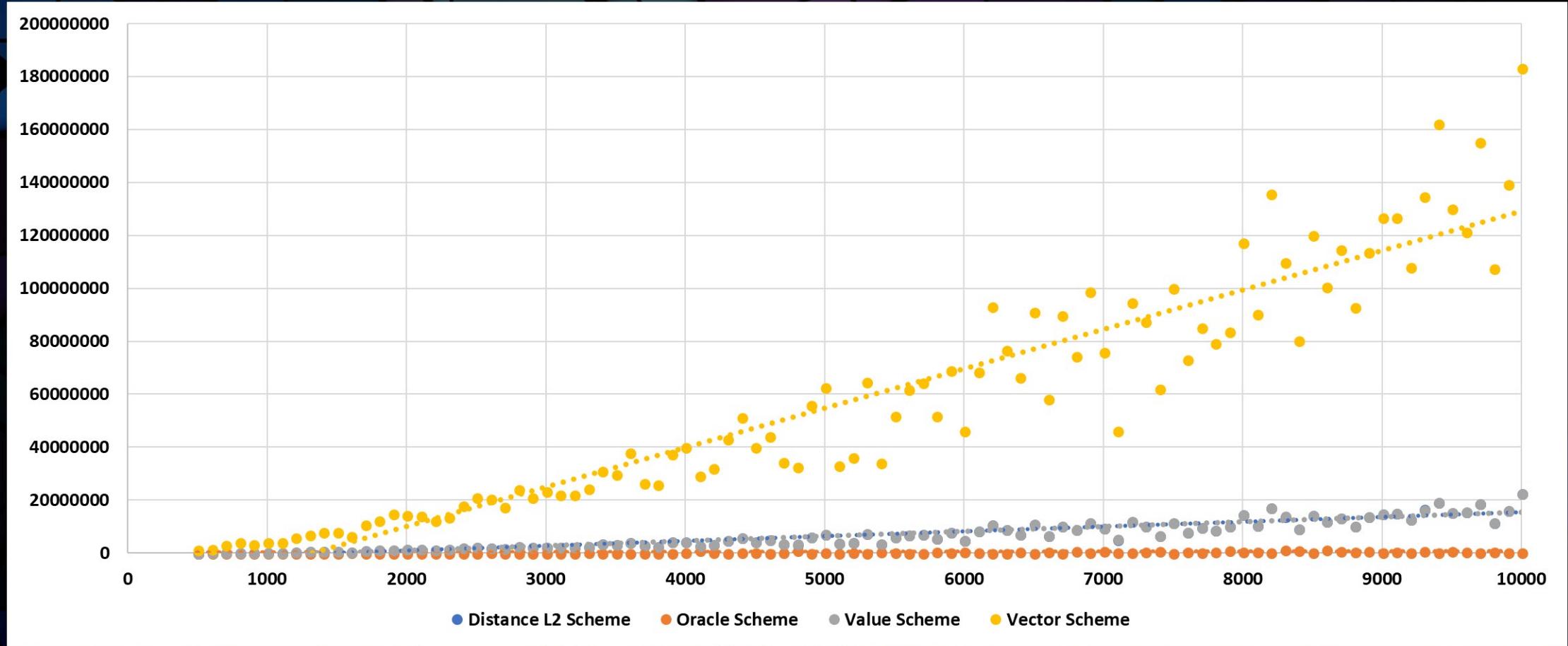
Global, Inner Data, $|\text{vector}| = 10000$

Time



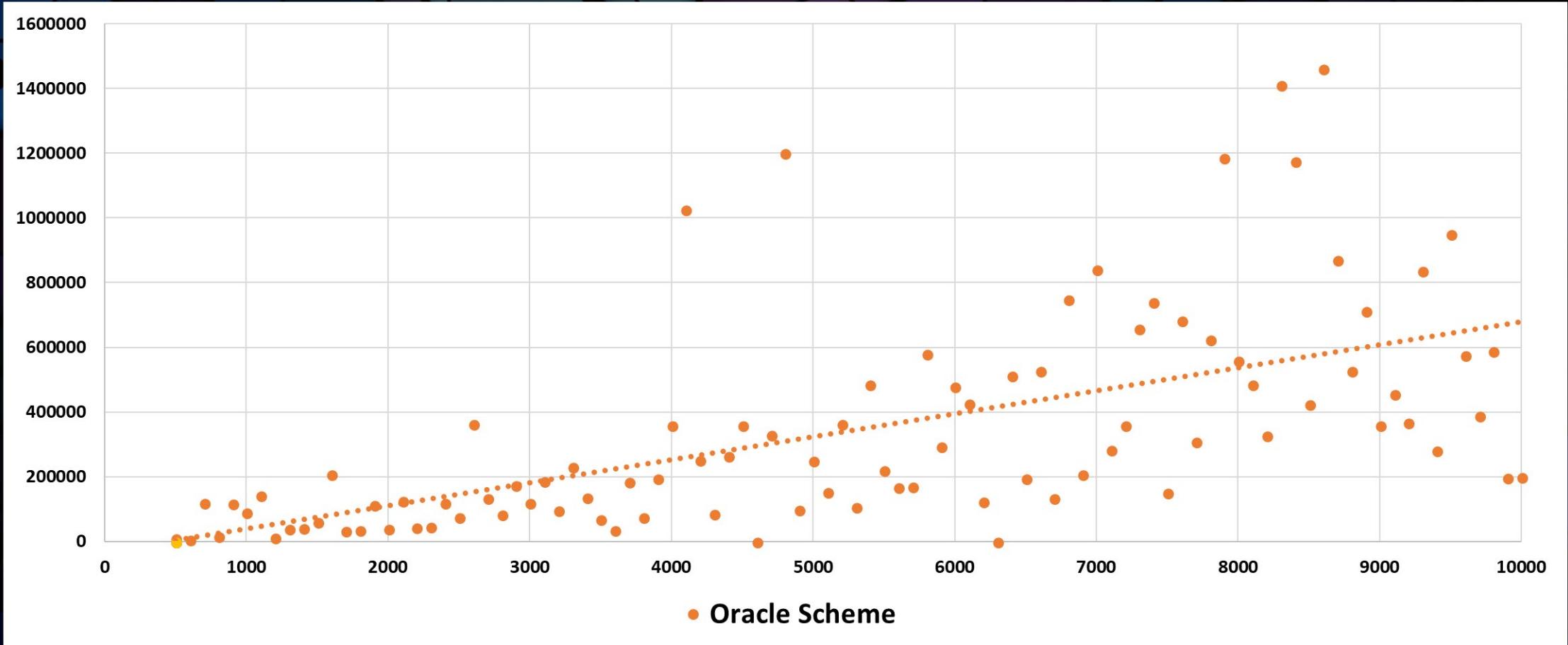
Bandwidth / $|\text{vector}|$

Bandwidth
↑
 $\rightarrow |\text{vector}|$



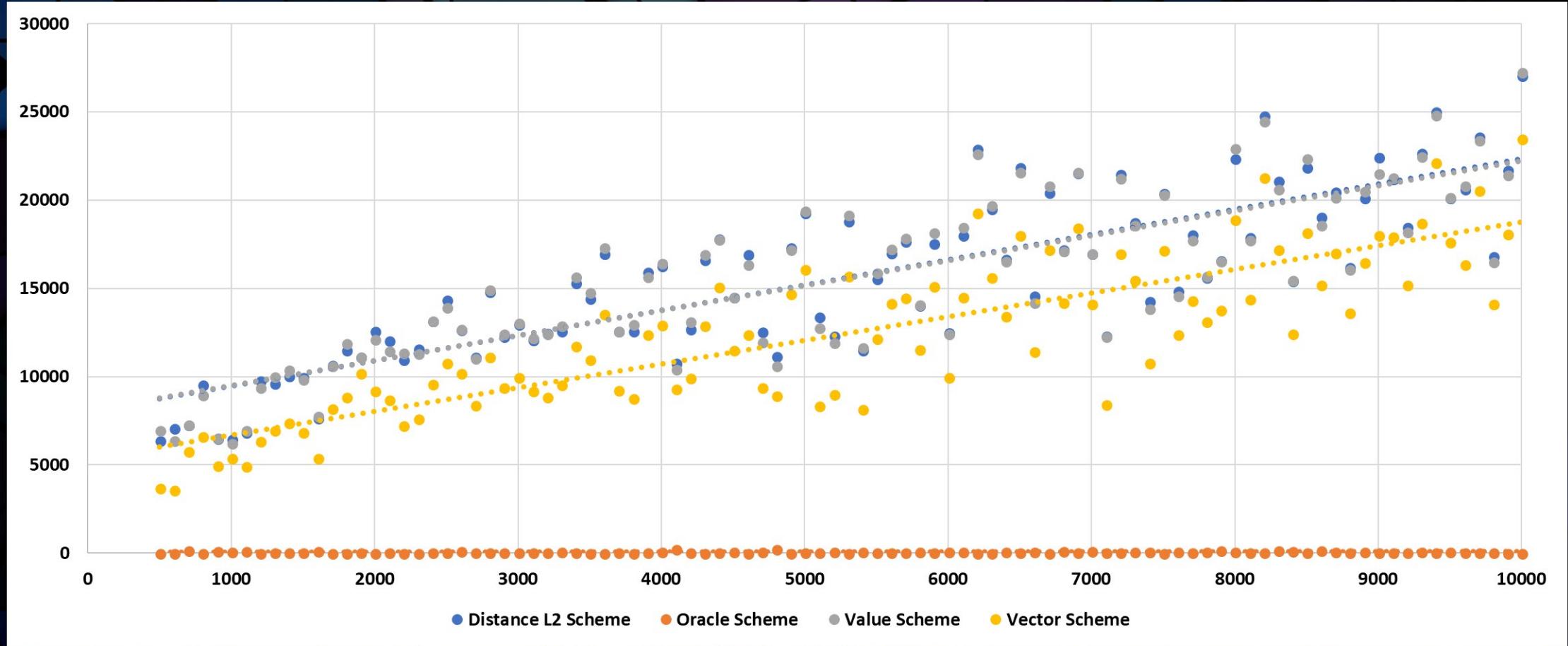
Bandwidth / $|\text{vector}|$

Bandwidth
 \uparrow
 $\rightarrow |\text{vector}|$



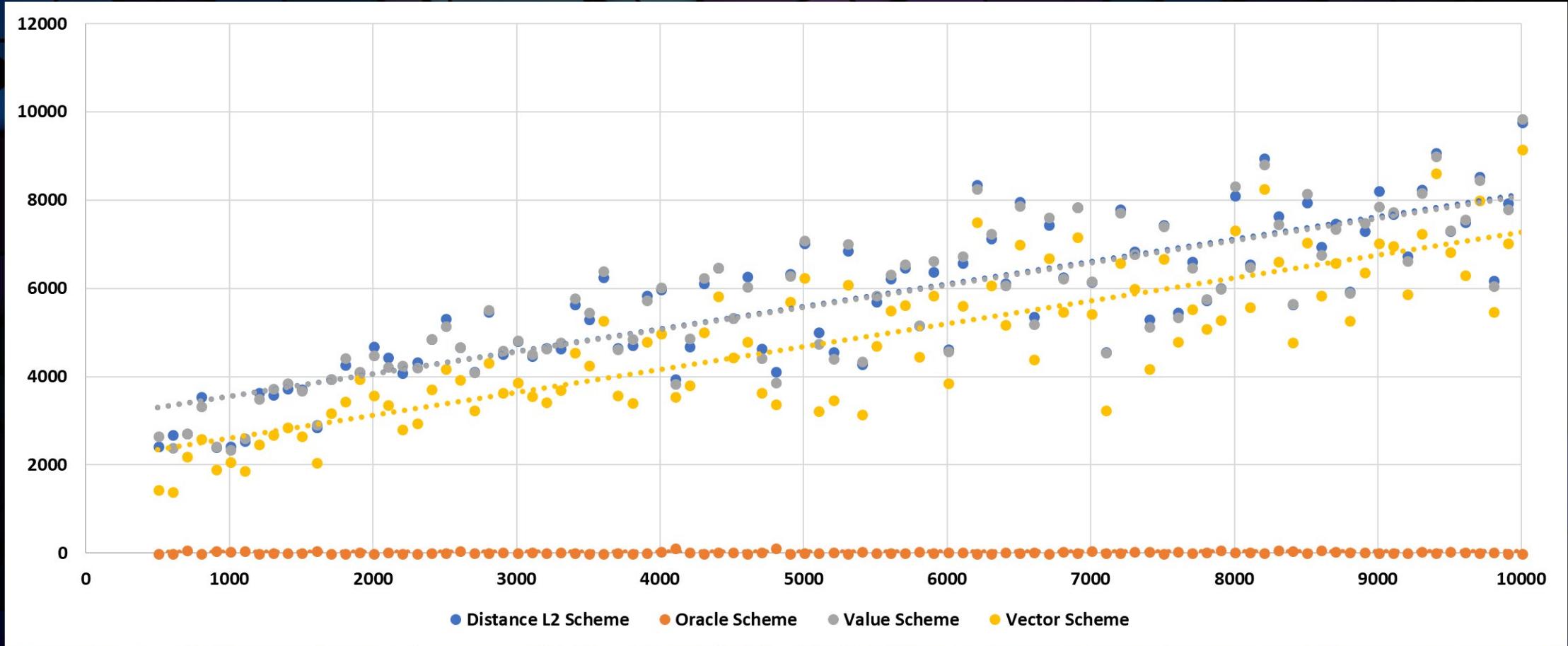
Messages / |vector|

Messages
↑
→ |vector|



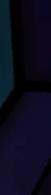
Channels / |vector|

Channels
↑
→ |vector|

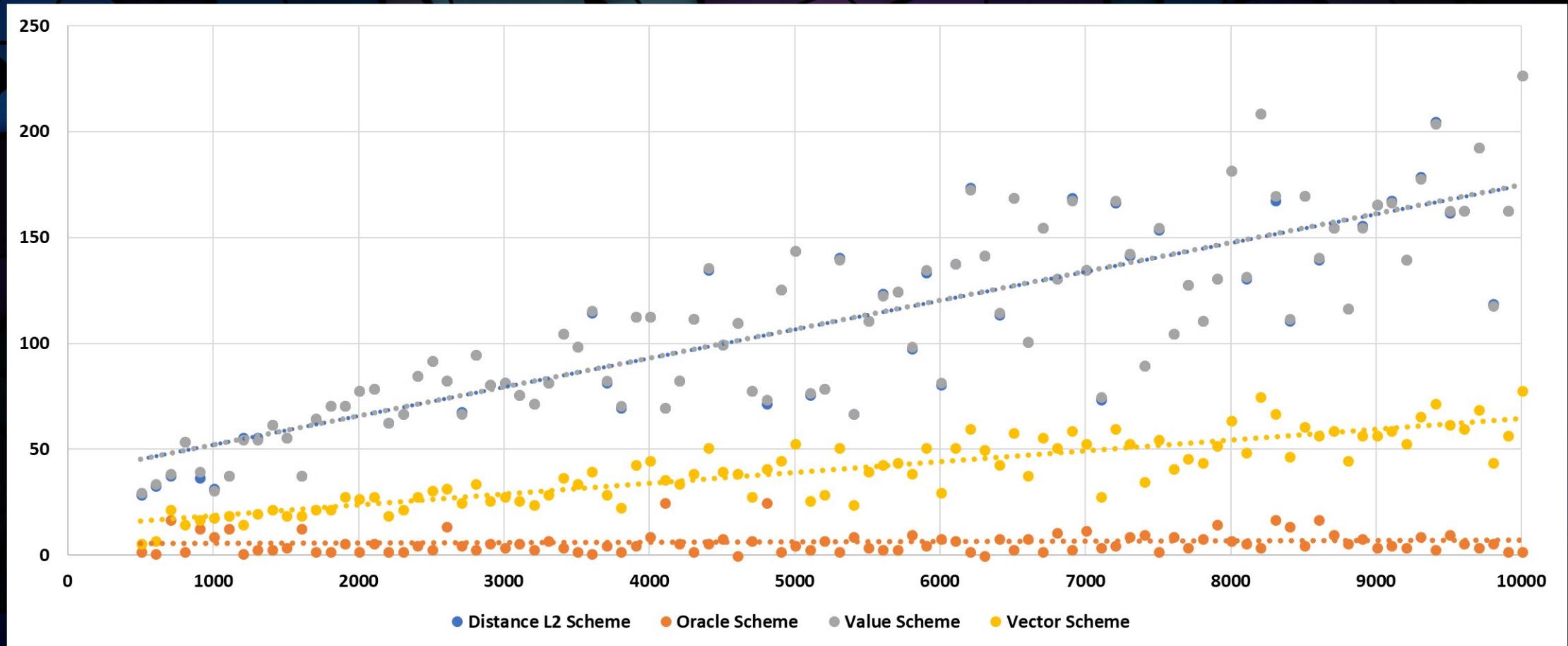


full syncs / $|\text{vector}|$

Full Syncs

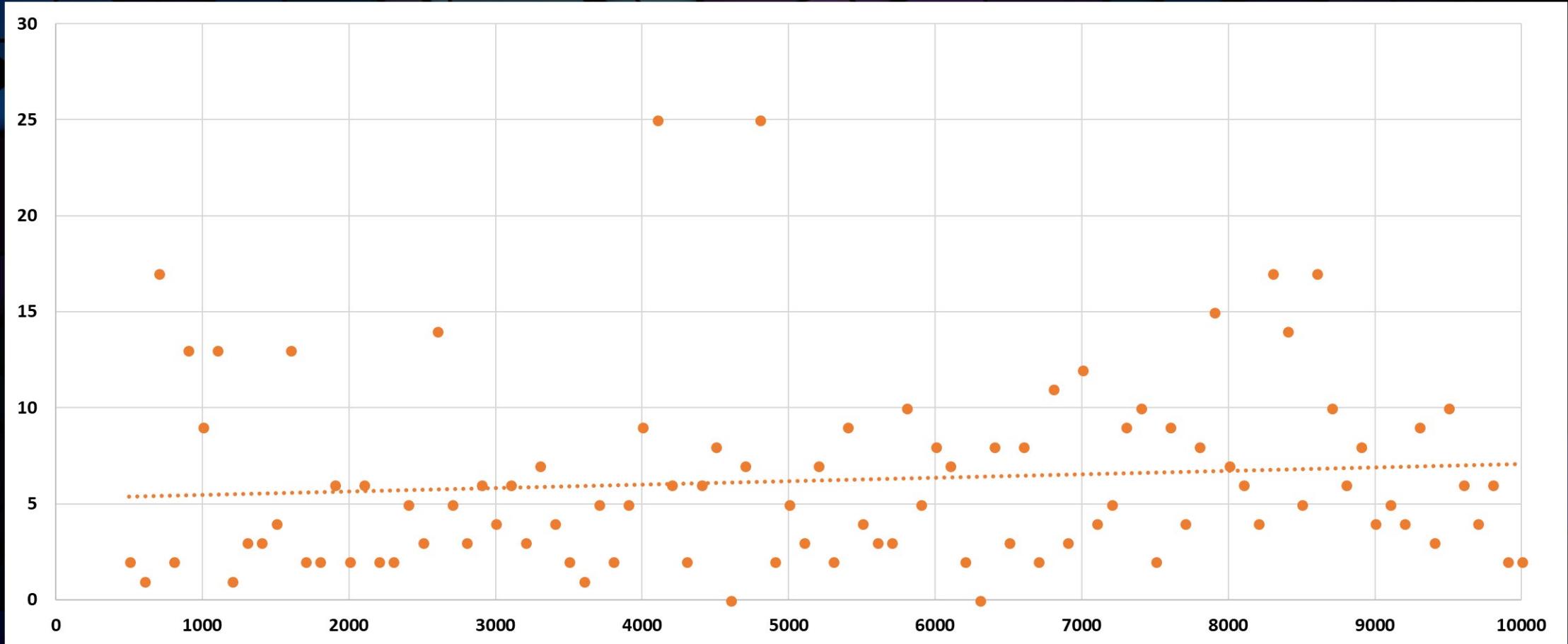


$|\text{vector}|$



ORACLE # full syncs / $|\text{vector}|$

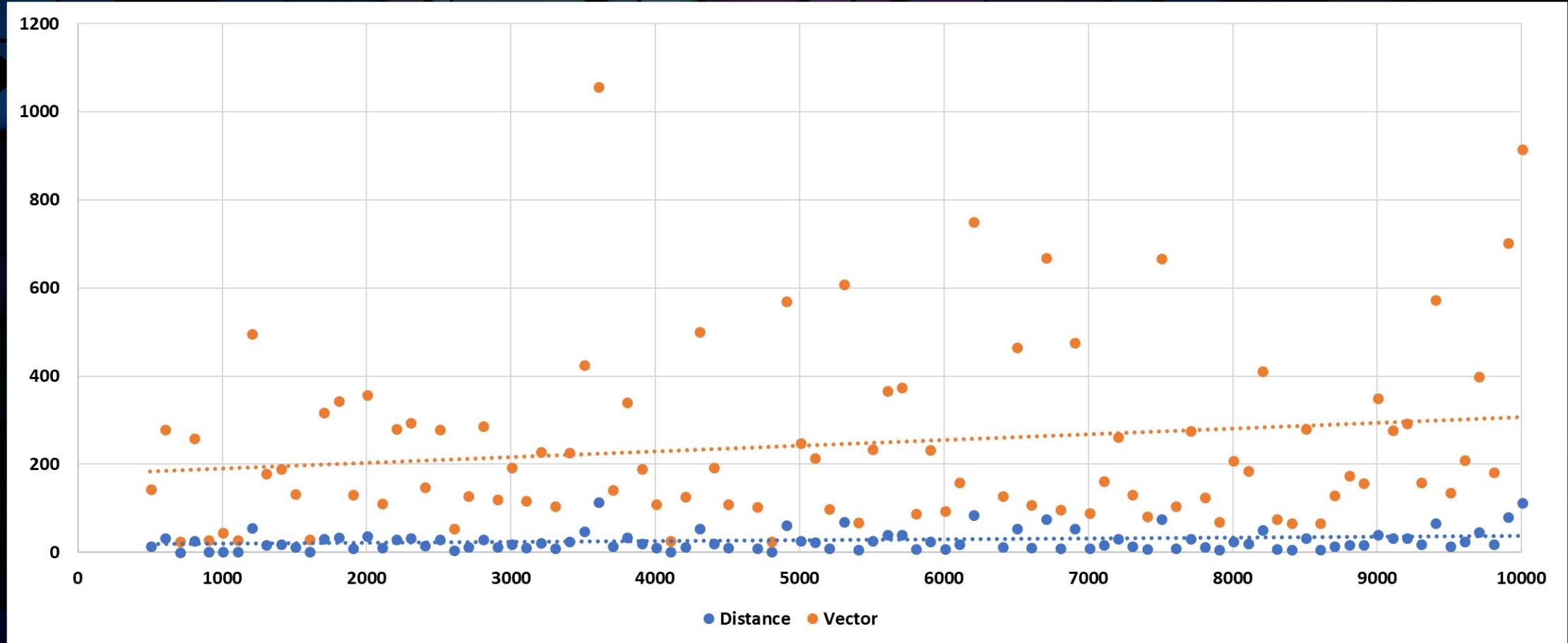
Full Syncs
↑
→ $|\text{vector}|$



$\frac{\text{Scheme Bandwidth}}{\text{Oracle Bandwidth}} / |\text{vector}|$

Scheme Bandwidth

Oracle Bandwidth



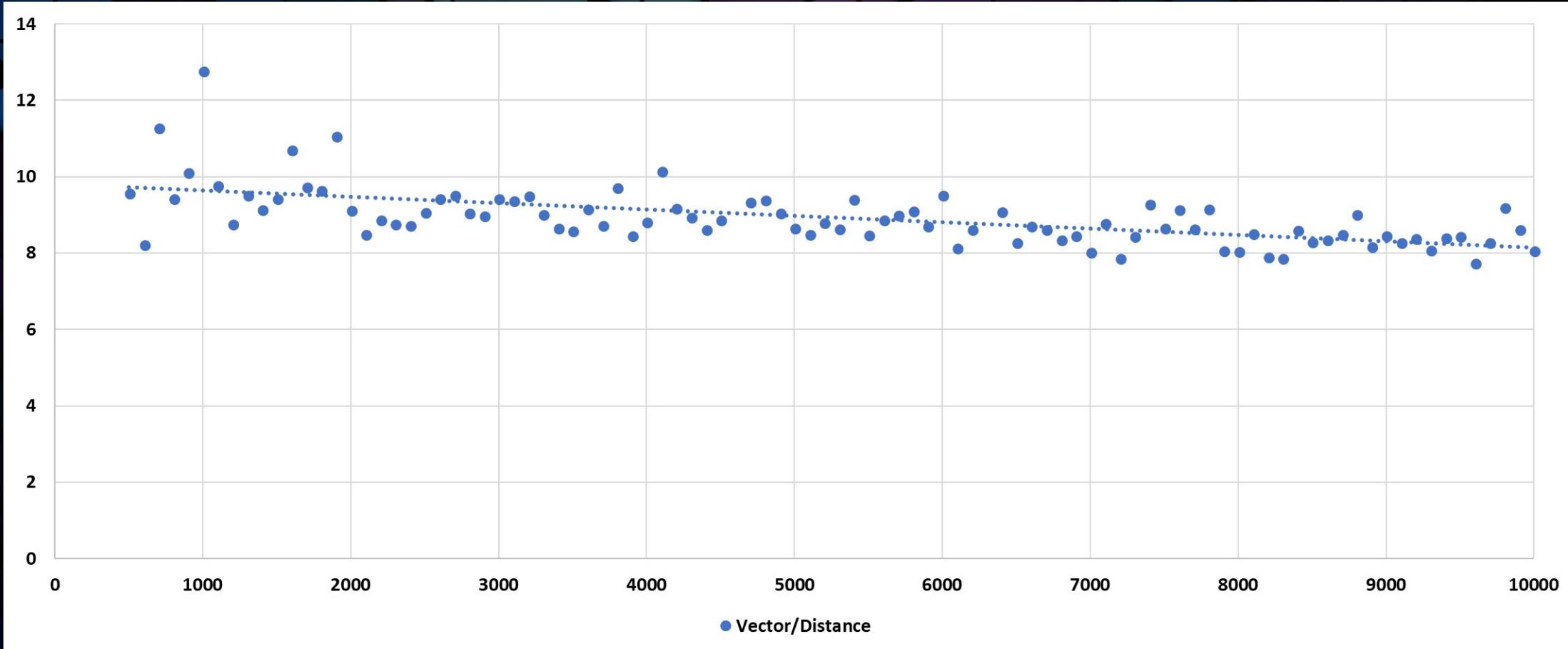
$\frac{\text{Vector Bandwidth}}{\text{Distance Bandwidth}} / |\text{vector}|$

Vector Bandwidth

Distance Bandwidth



$|\text{vector}|$





Entropy

Entropy Monitoring

- ❑ Distance Scheme – L_1 Distance

- ❑ Data Vector:

`vector(i) <- # occurrences in the window of the i'th lexicographic ordered n-most common token`

$$\text{Entropy}(\text{vector}) = \sum p_i \cdot \ln(1/p_i)$$

Running Information

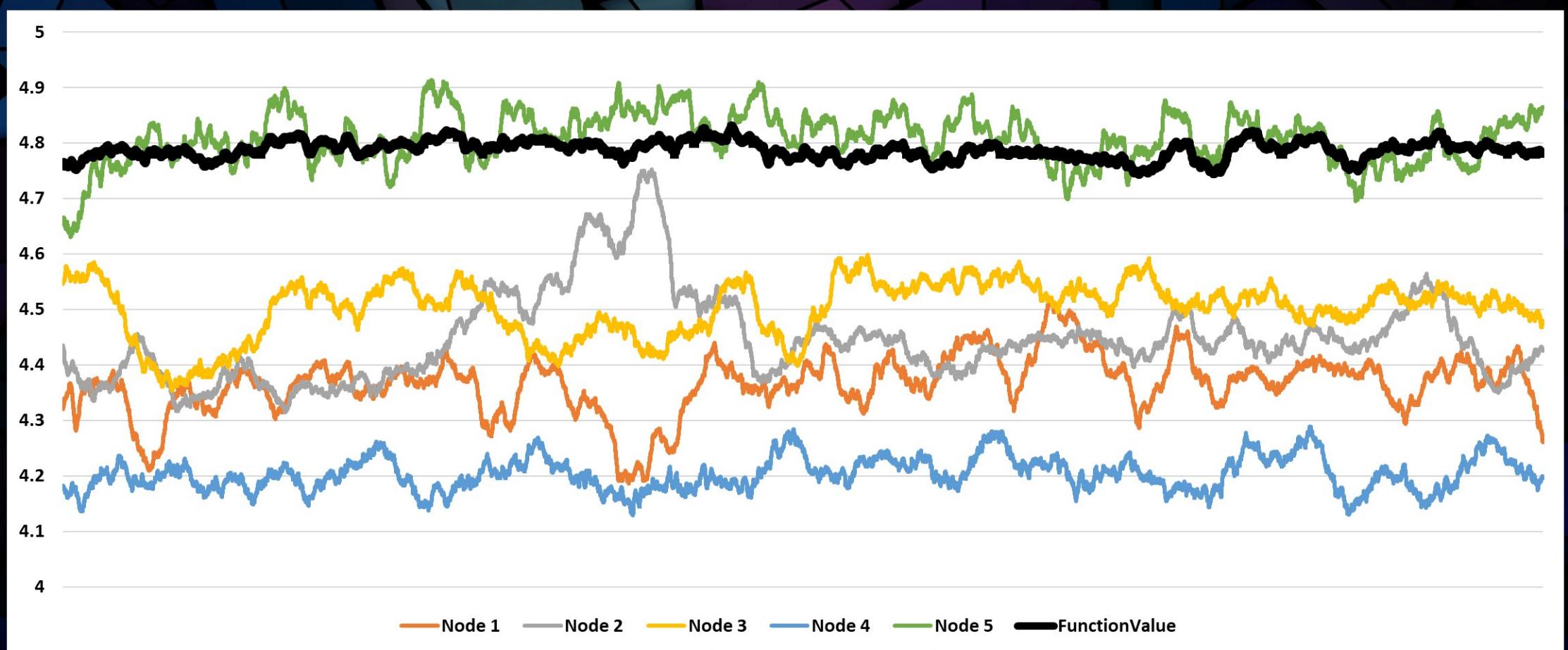
- Window Size = 5000
- Step Size = 20
- Epsilon – Multiplicative (*0.995, *1.005)

Entropy



Time

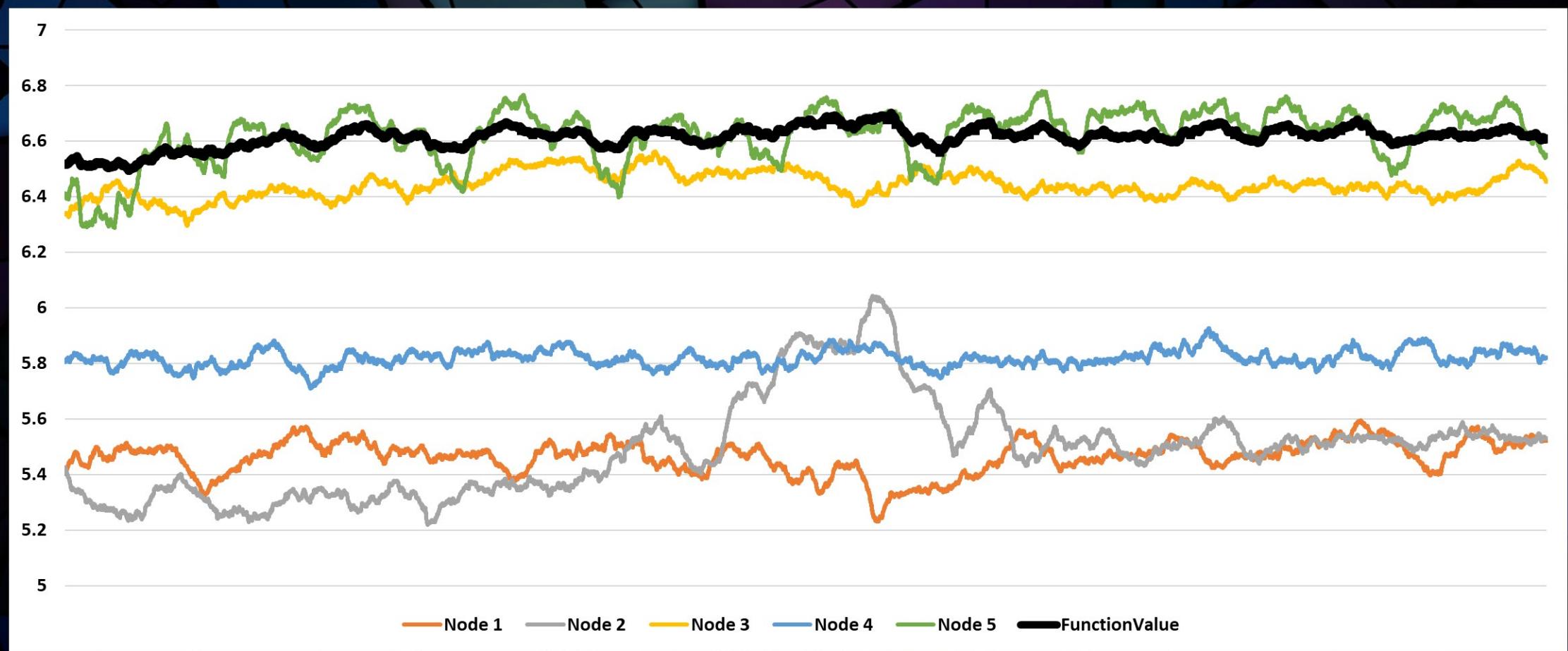
Entropy Data, $|\text{vector}| = 500$



Entropy

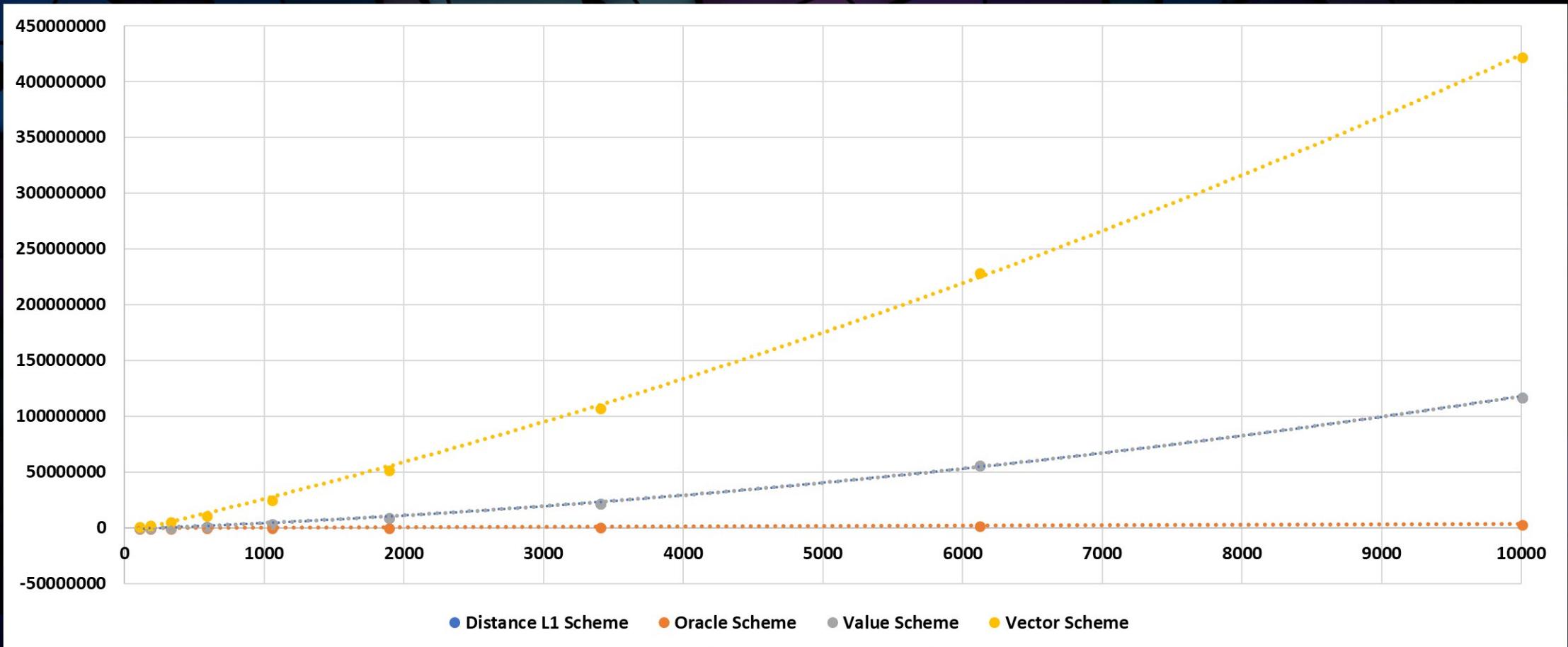
Entropy Data, $|\text{vector}| = 10000$

Time



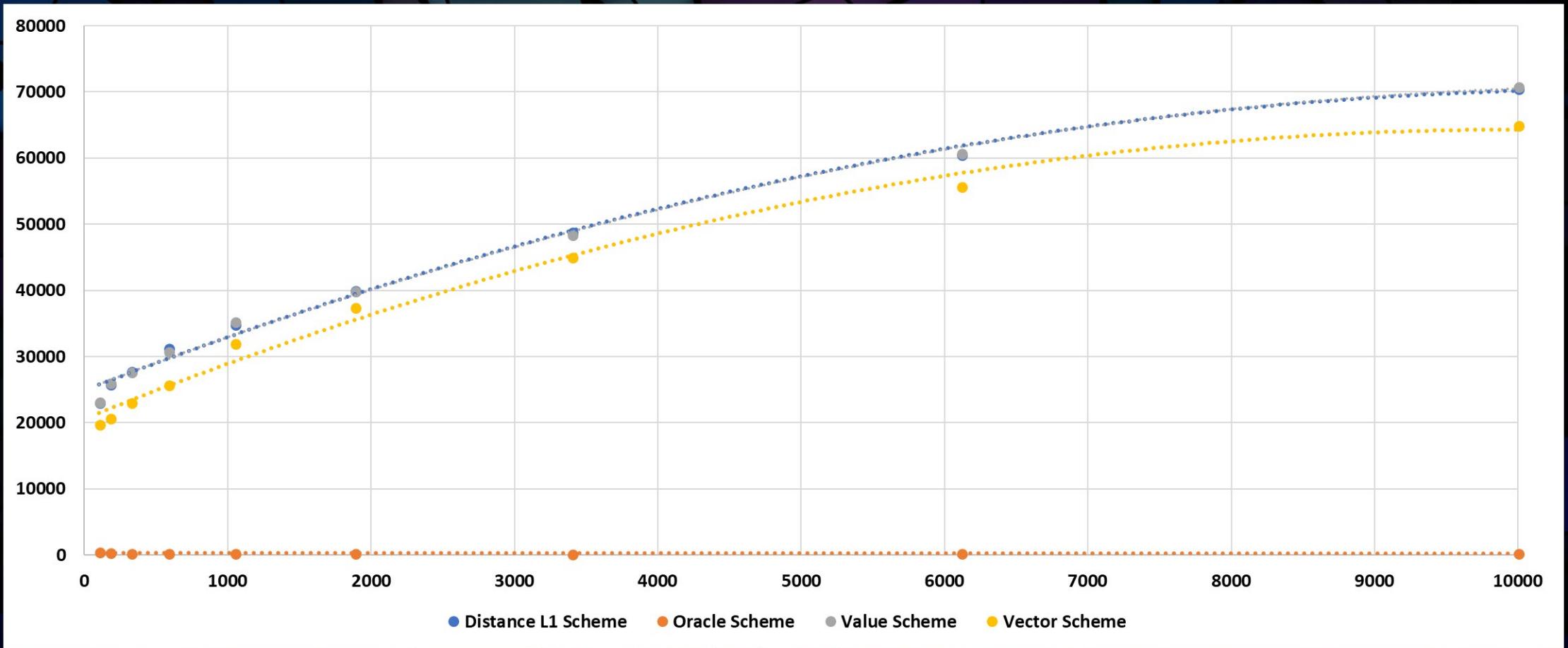
Bandwidth / $|\text{vector}|$

Bandwidth
 \uparrow
 $\rightarrow |\text{vector}|$



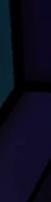
Messages / |vector|

Messages
↑
→ |vector|

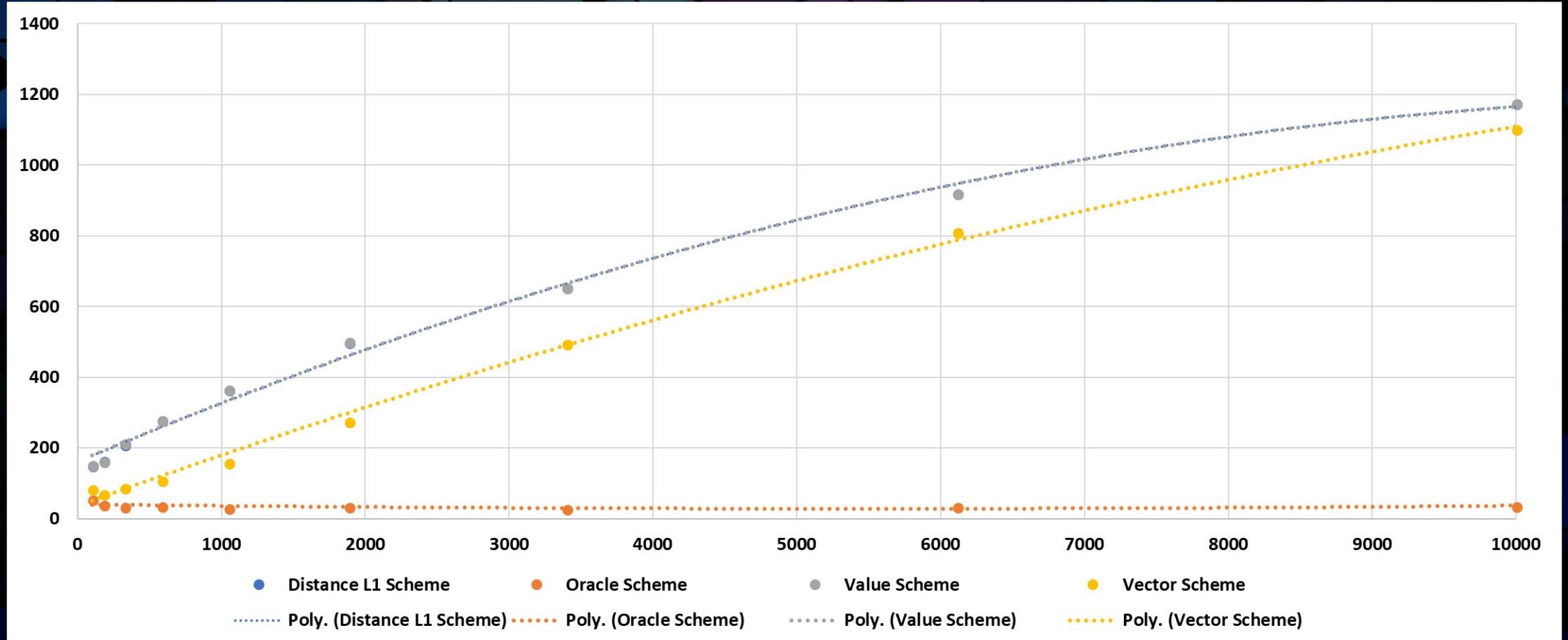


full syncs / $|\text{vector}|$

Full Syncs



$|\text{vector}|$

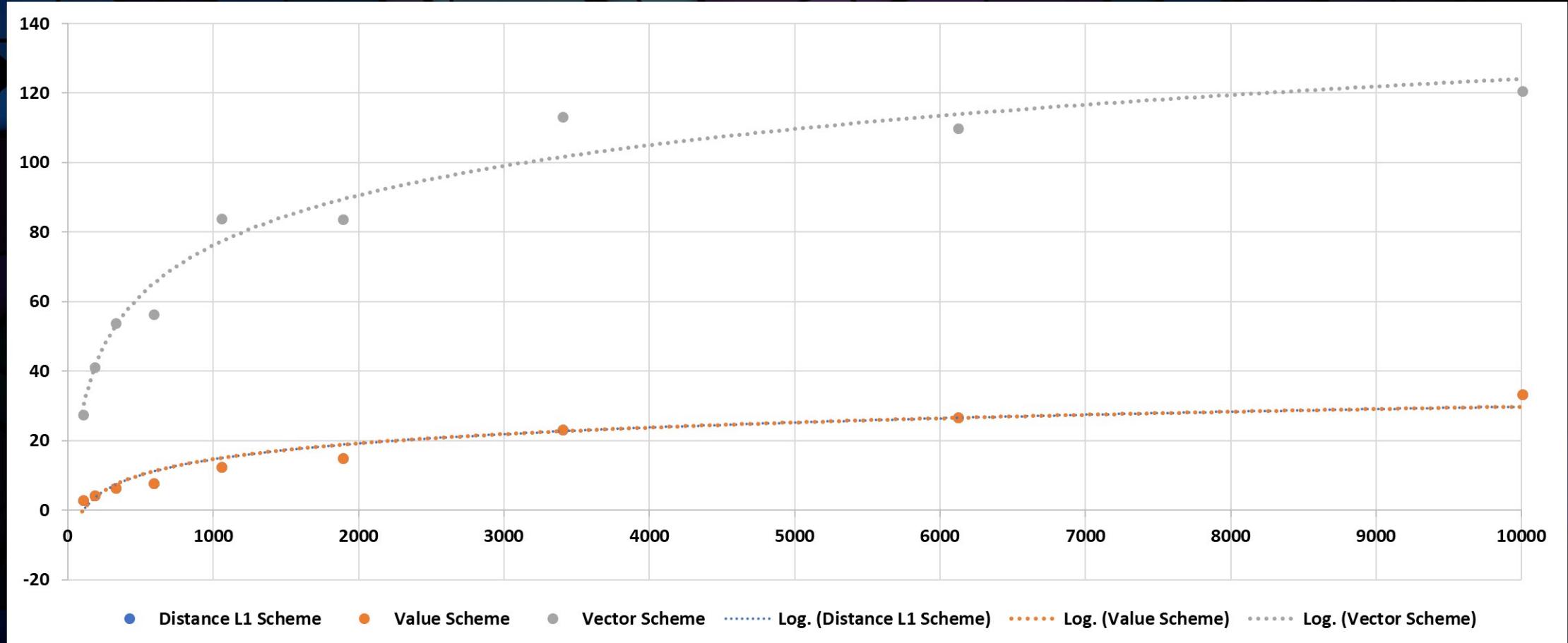


$\frac{\text{Scheme Bandwidth}}{\text{Oracle Bandwidth}}$

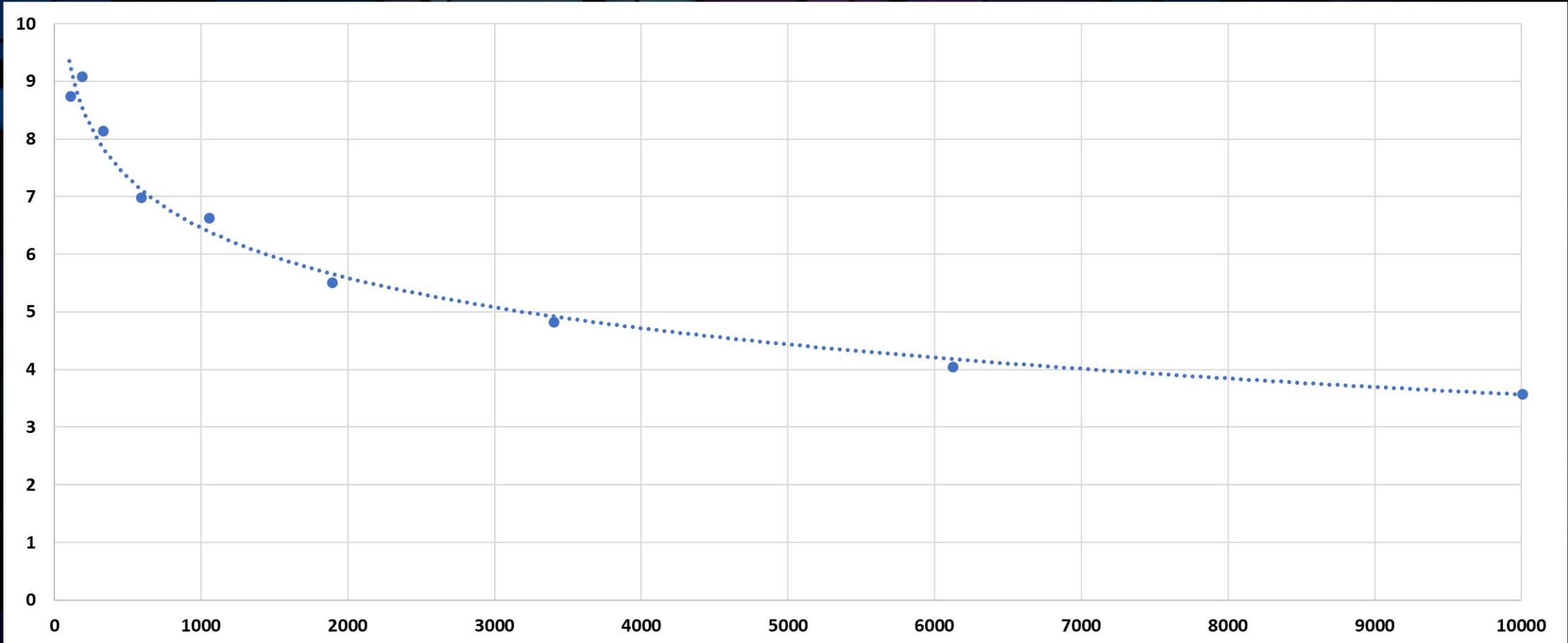
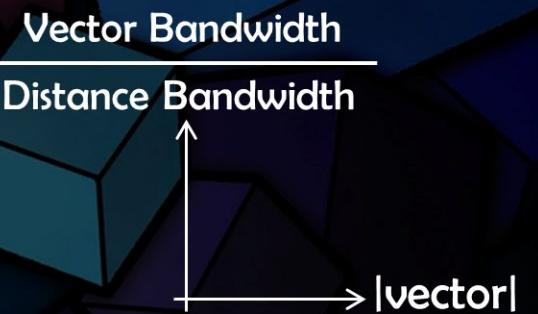
/ |vector|

Scheme Bandwidth

Oracle Bandwidth



$\frac{\text{Vector Bandwidth}}{\text{Distance Bandwidth}} / |\text{vector}|$





PCA

PCA

- PCA of residual vectors could further decrease the bandwidth in Distance Scheme
- In order to asses how could will it perform, we should look at the eigenvalues of the PCA of all the residual vectors
- Problem:

PCA takes cubic time in the dimensionality of the data

$|\text{vector}| = 440 \rightarrow 2 \text{ minute}$

$|\text{vector}| = 1000 \rightarrow 24 \text{ minutes}$

$|\text{vector}| = 1700 \rightarrow 2 \text{ hour}$

$|\text{vector}| = 3900 \rightarrow 2 \text{ day}$

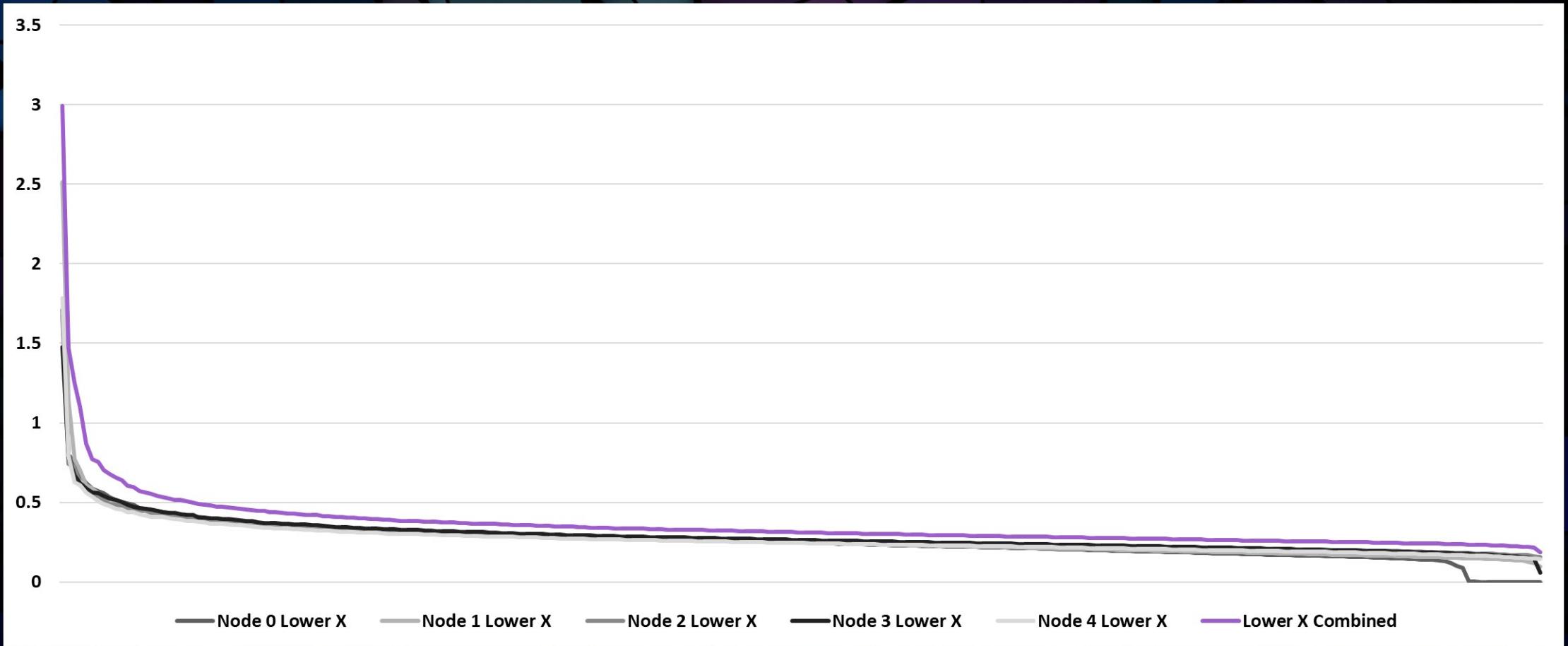
Inner Product PCA

- Inner Product
- Vector Length = 5,000
- Window Size = 50,000
- Step Size = 100
- # Nodes = 5

Eigenvalue's
value, 10th
root scale

Nodes Lower Bound X Residuals

→ # eigenvalue

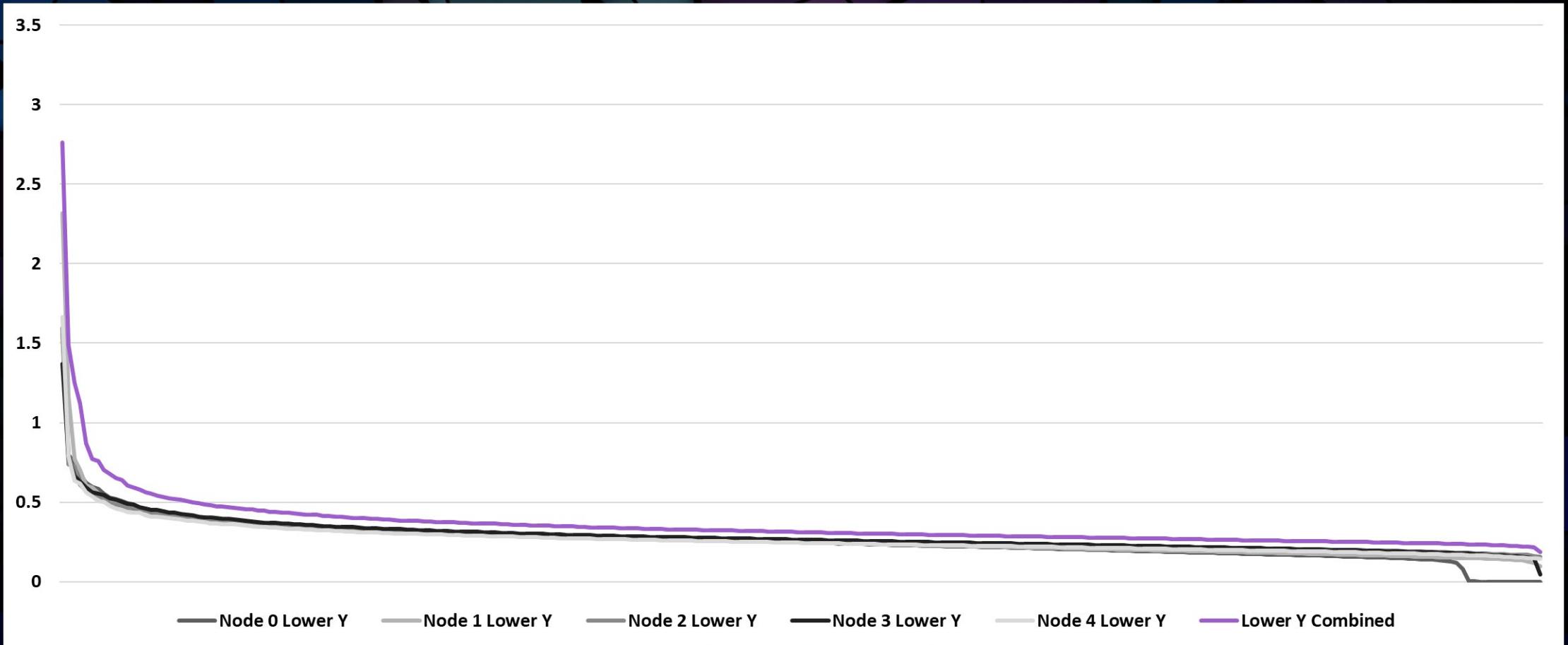


Eigenvalue's
value, 10th
root scale



→ # eigenvalue

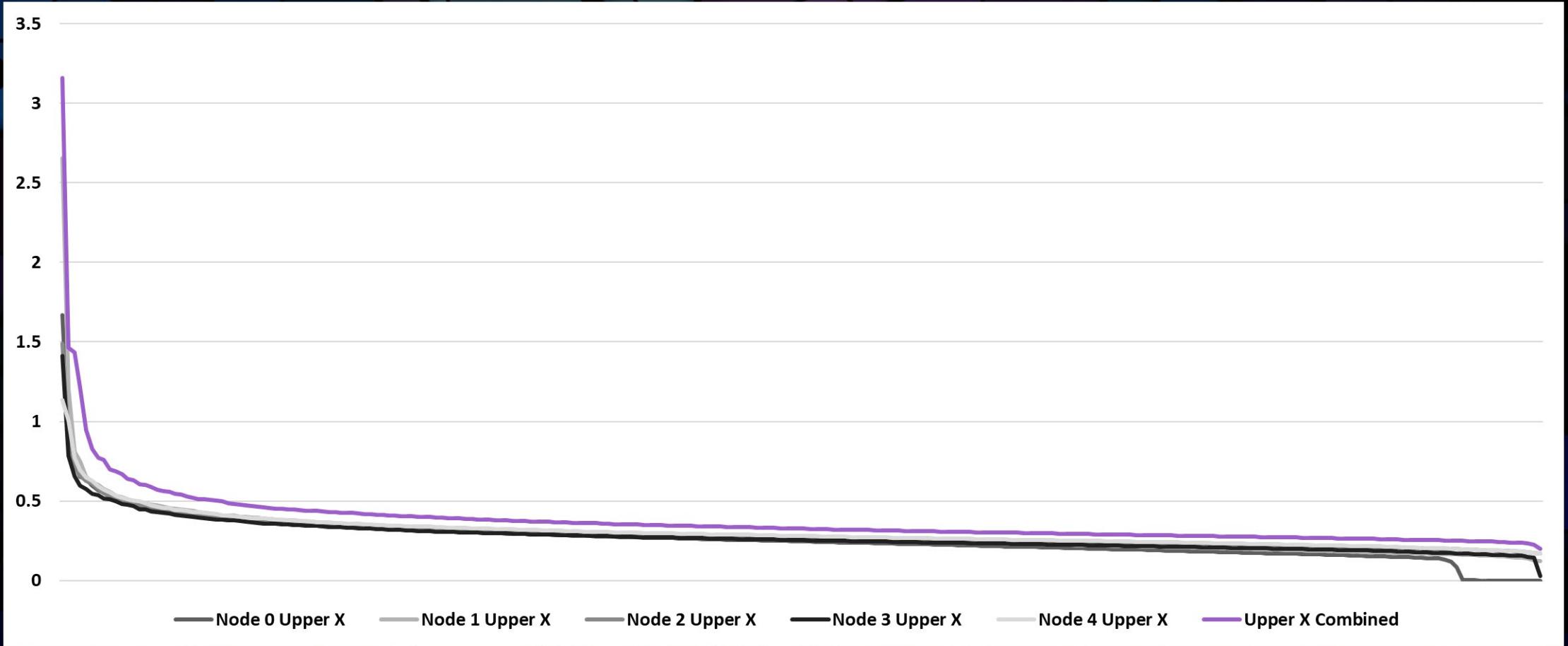
Nodes Lower Bound Y Residuals



Eigenvalue's
value, 10th
root scale

Nodes Upper Bound X Residuals

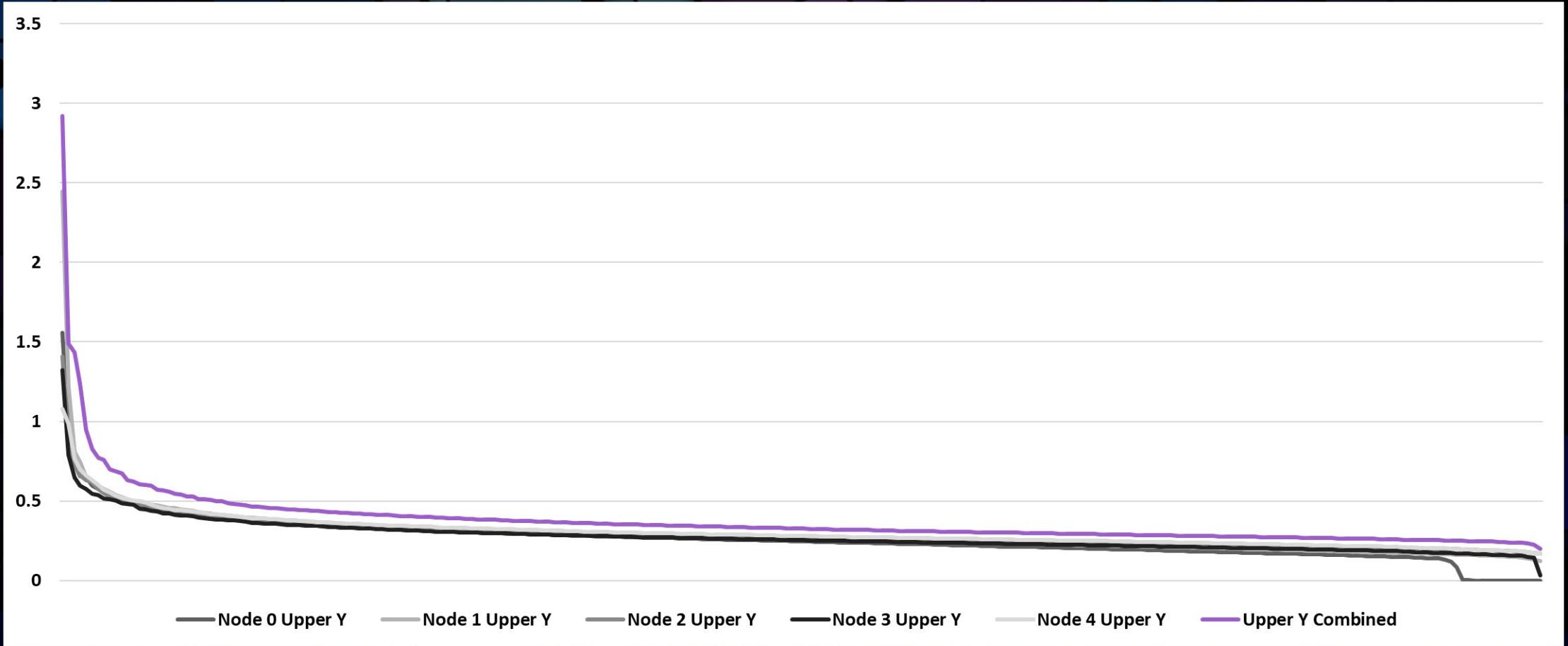
→ # eigenvalue



Eigenvalue's
value, 10th
root scale

Nodes Upper Bound Y Residuals

→ # eigenvalue



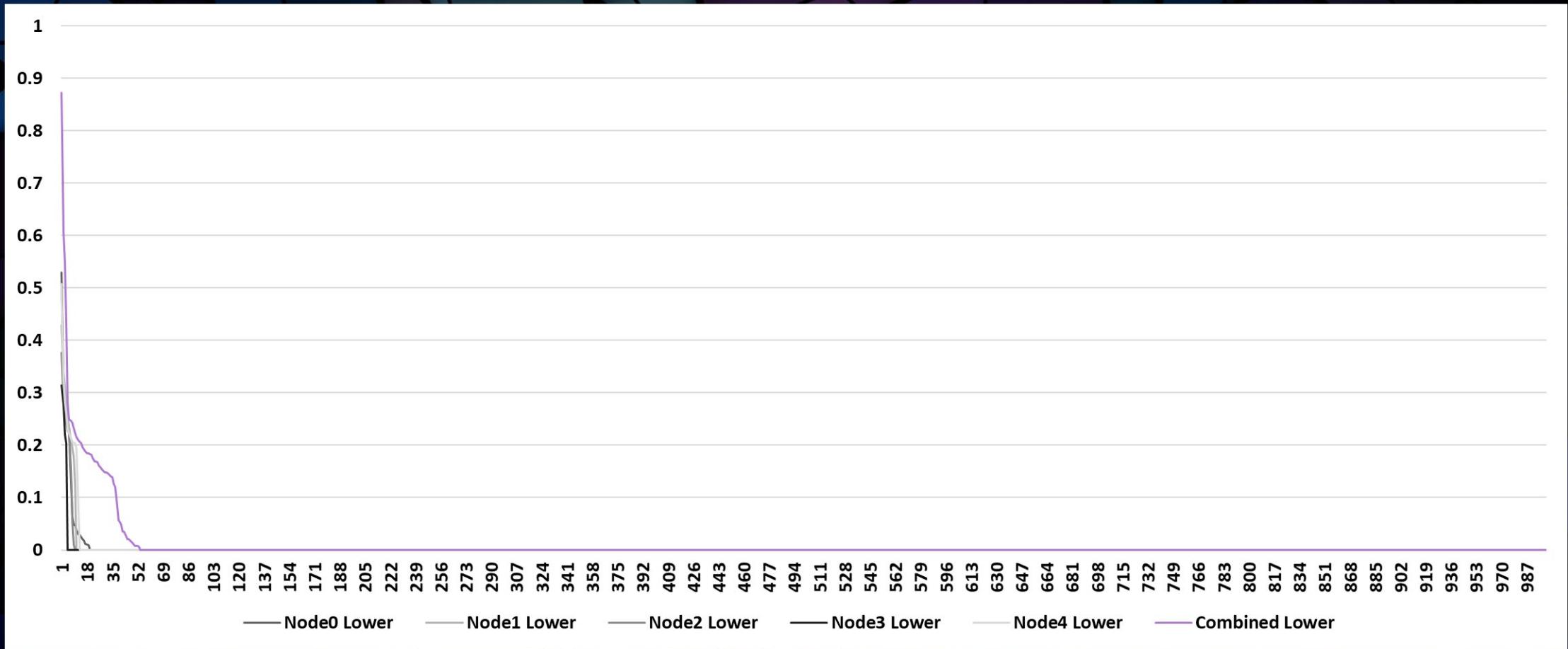
Entropy PCA

- Entropy
- Vector Length = 1,000
- Window Size = 10,000
- Step Size = 100
- # Nodes = 5

Eigenvalue's
value, 10th
root scale

→ # eigenvalue

Nodes Lower Bound Residuals



Eigenvalue's
value, 10th
root scale

Nodes Upper Bound Residuals

→ # eigenvalue

