



Exploring Cell Biology Literature with NLP

Beth Baumann

PubMed

 **GENSIM**
topic modelling for humans

spaCy



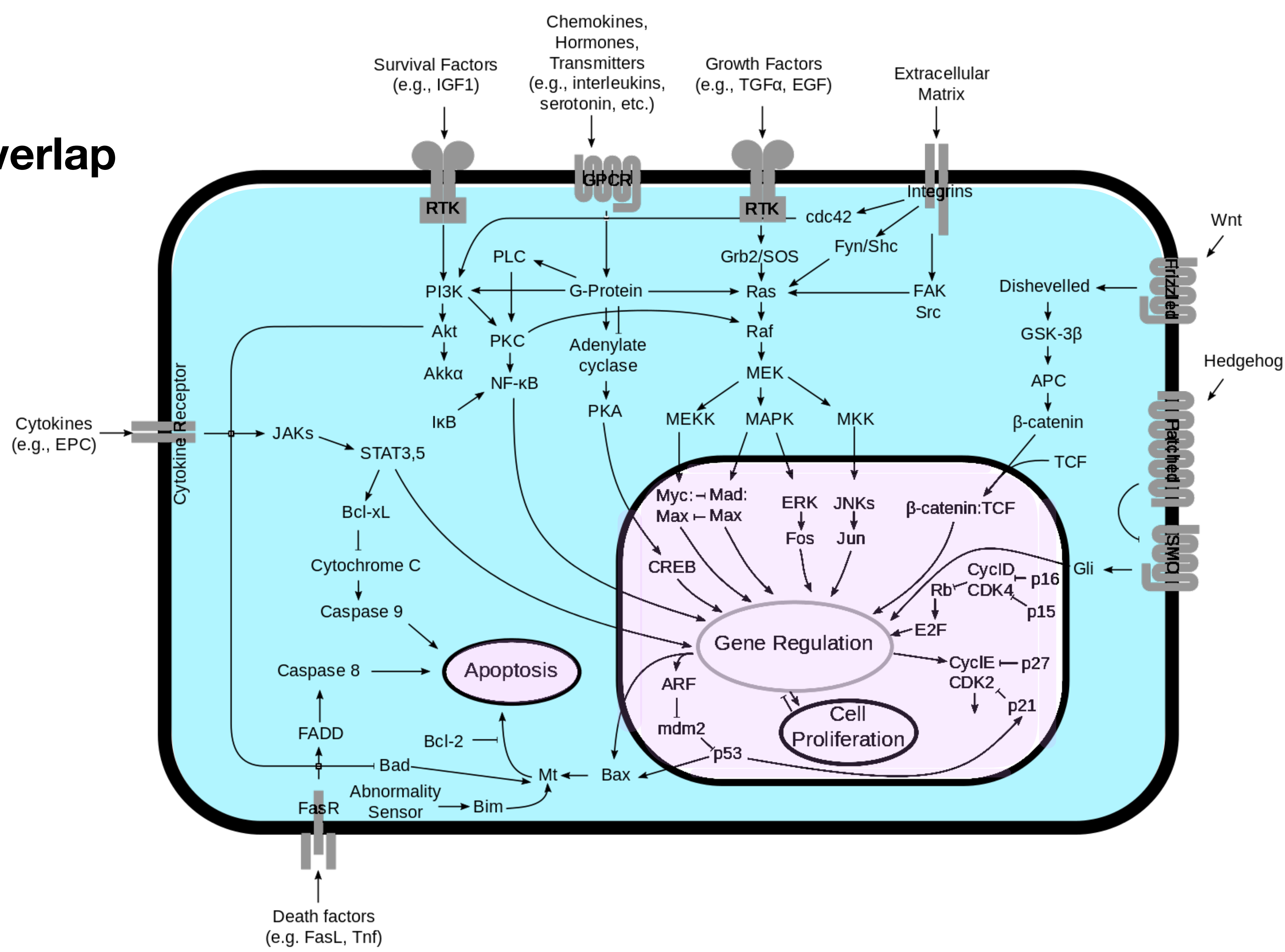
Goals of this NLP analysis

1. Identify the natural clustering of cell biology abstracts into signaling pathways

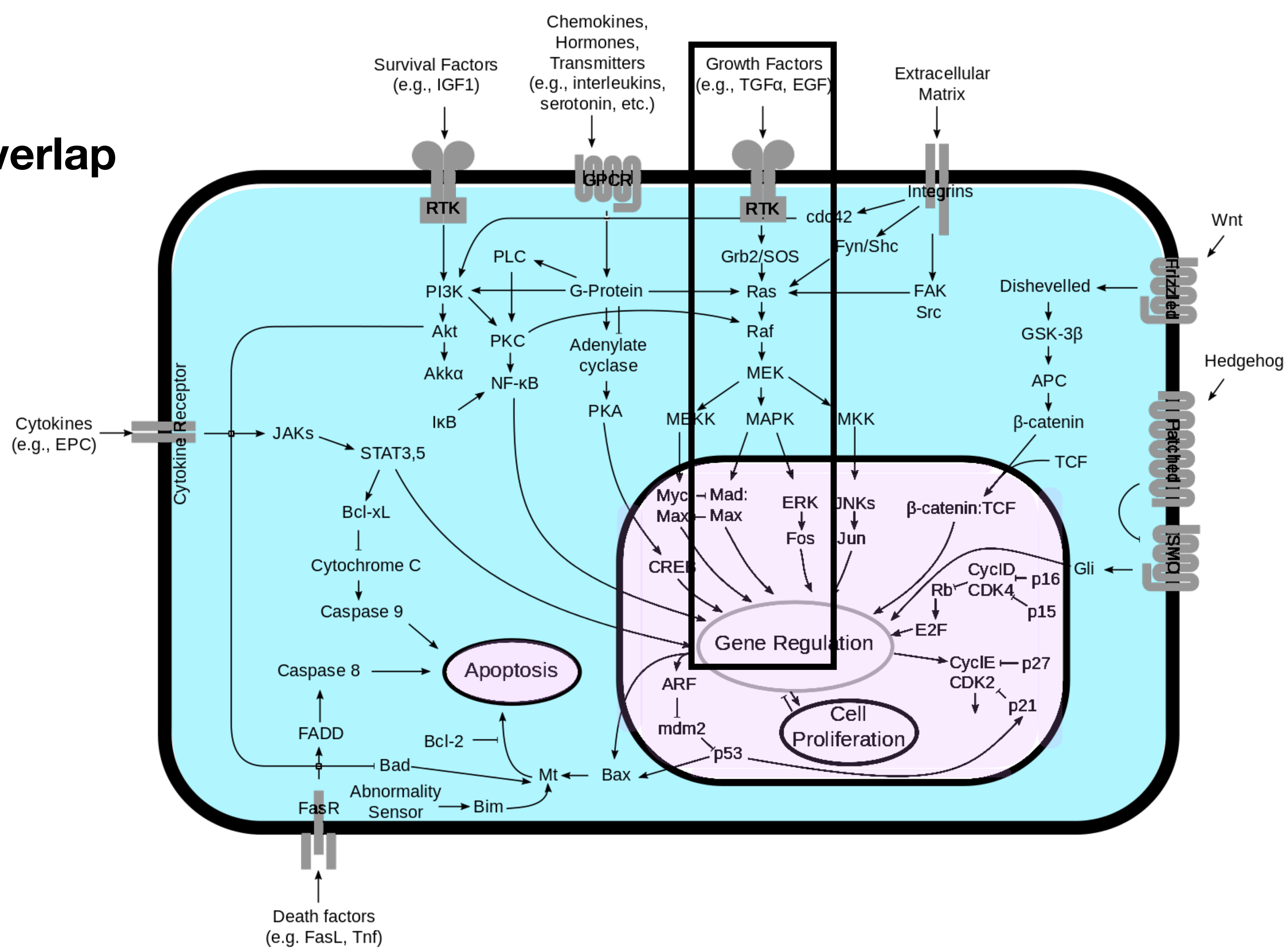
Signaling pathway: a chain or network of proteins that work together for a specific effect on the cell

2. Determine protein membership of signaling pathways
3. Explore cosine similarity of biological terms with Word2Vec

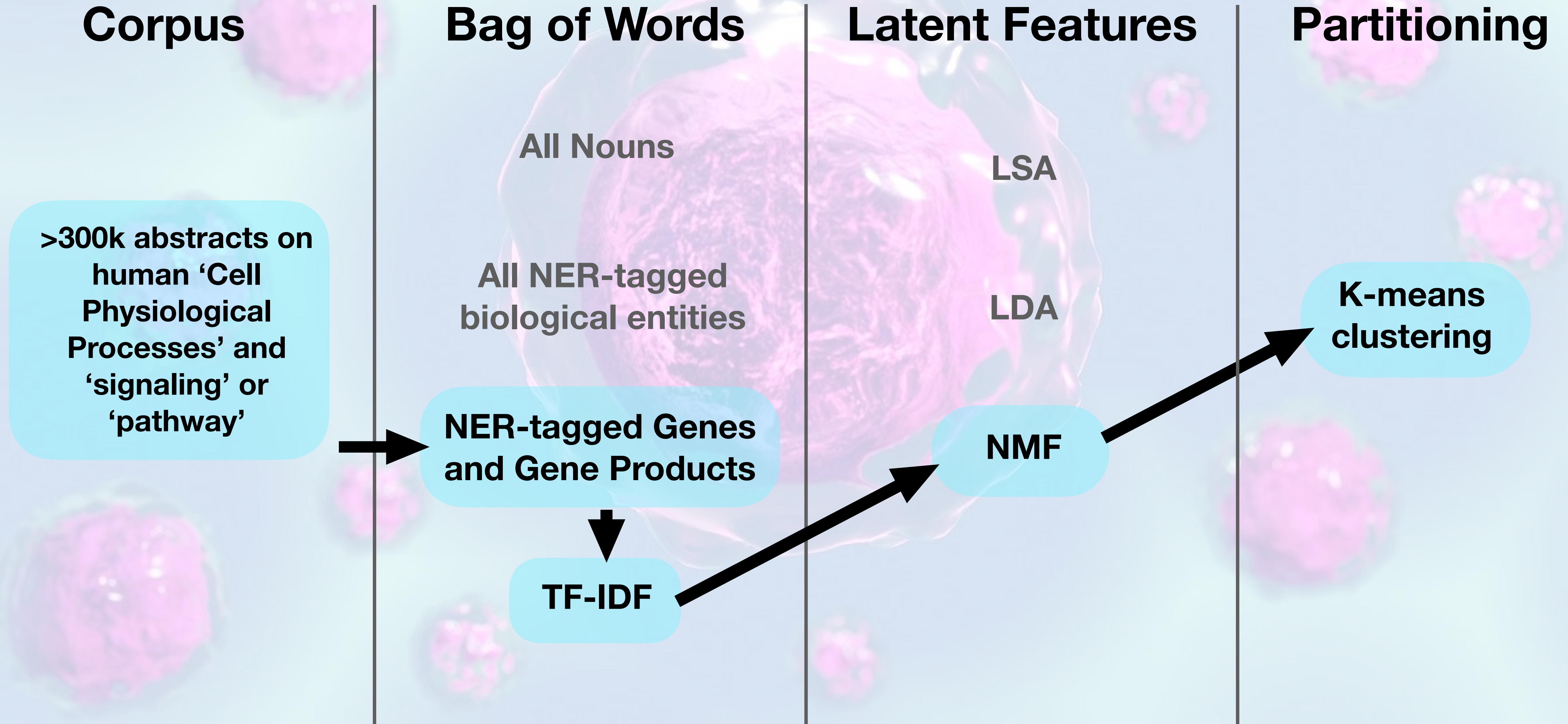
Signaling Pathway overlap



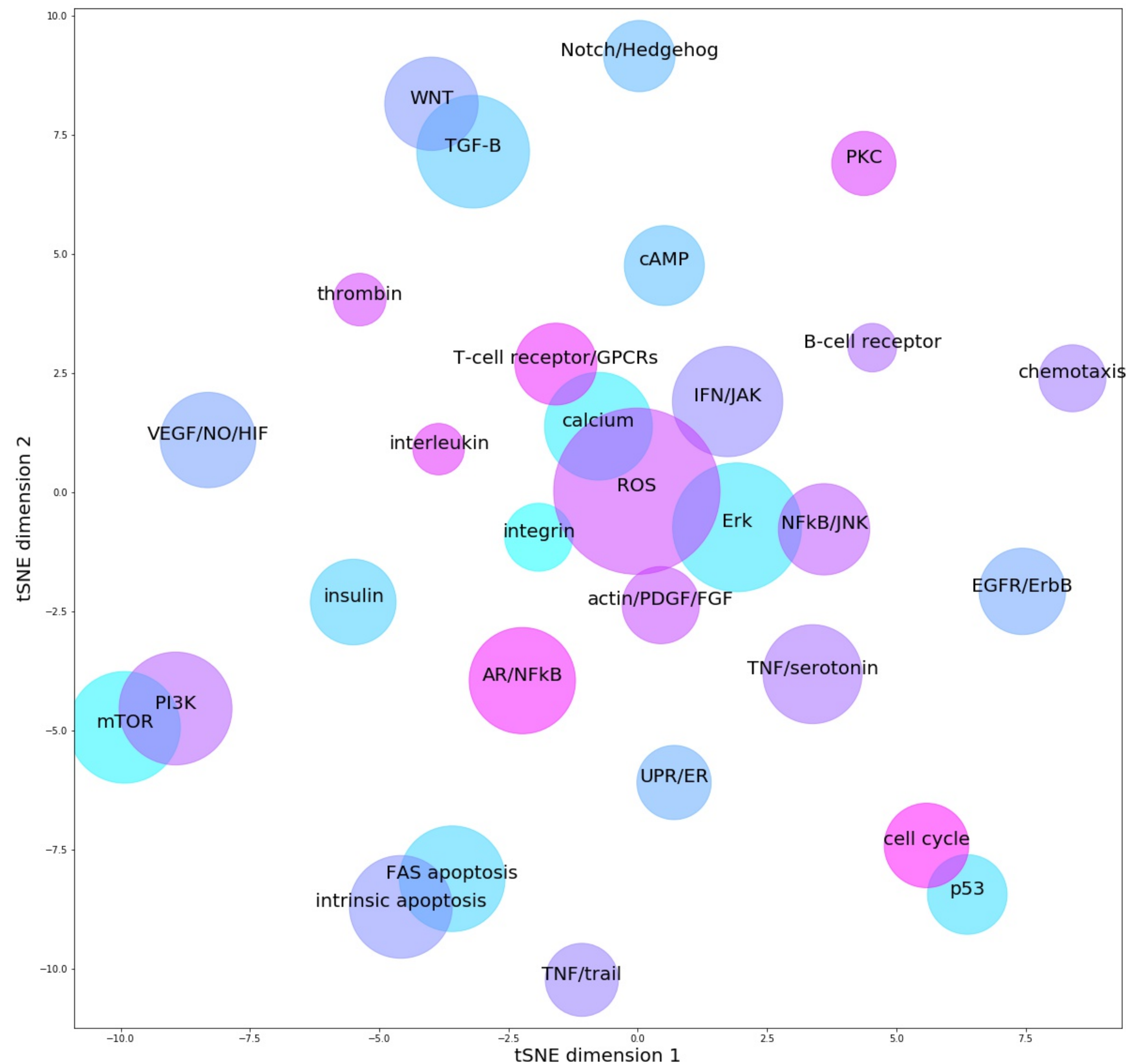
Signaling Pathway overlap



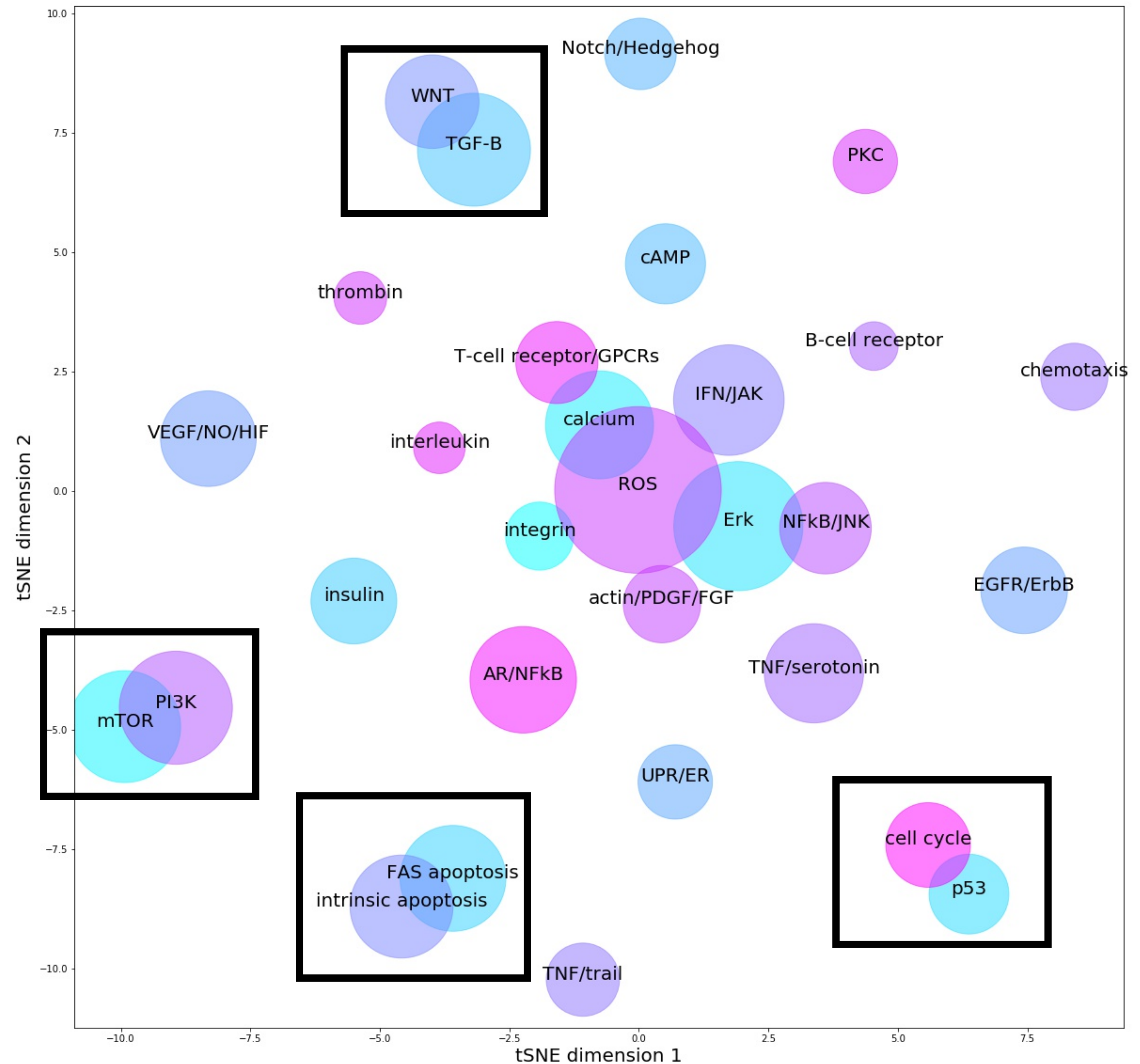
Part I: Topic Modeling with NMF



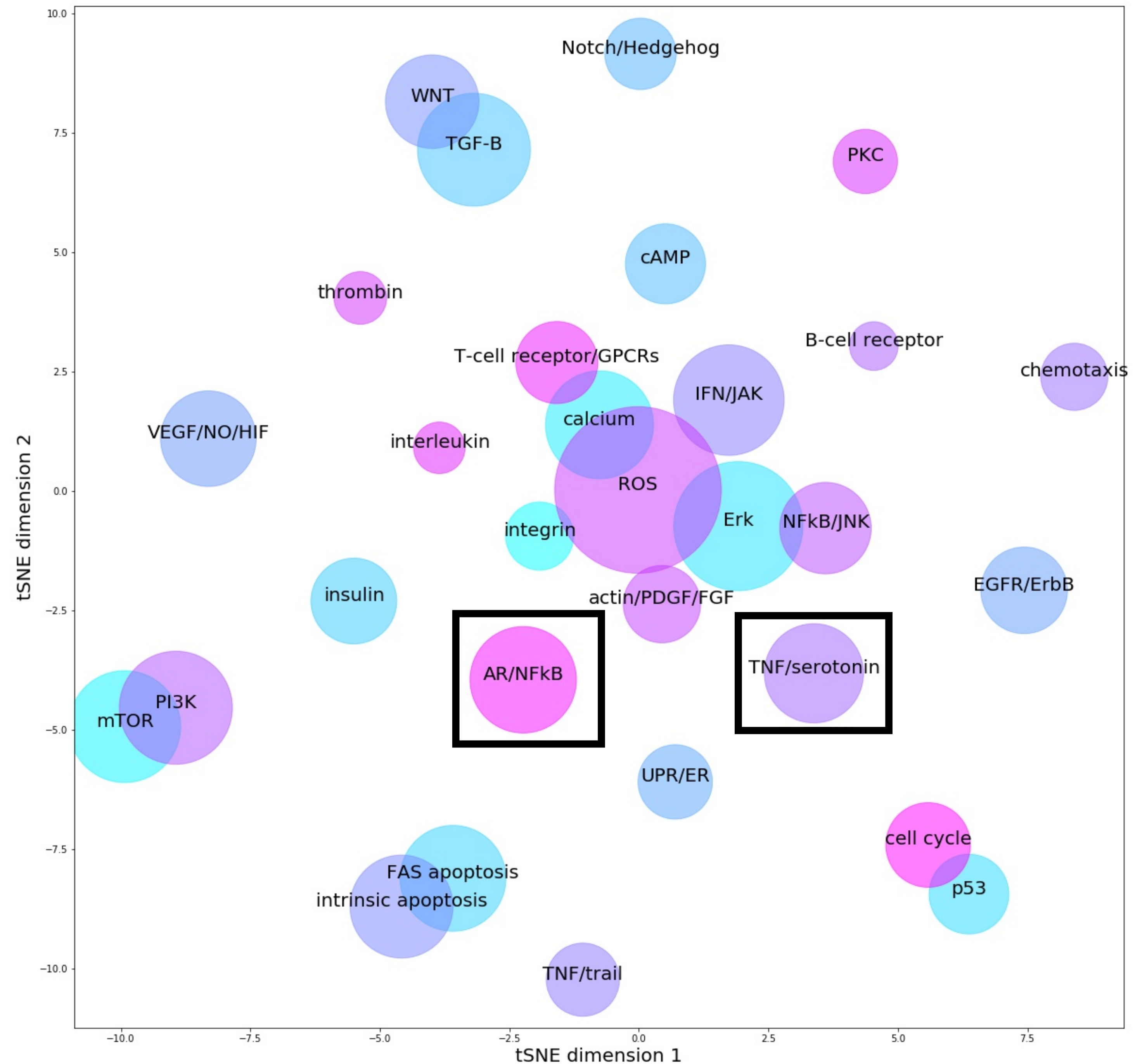
**Abstracts fell into
30 K-means clusters
along 25 latent topic
dimensions**



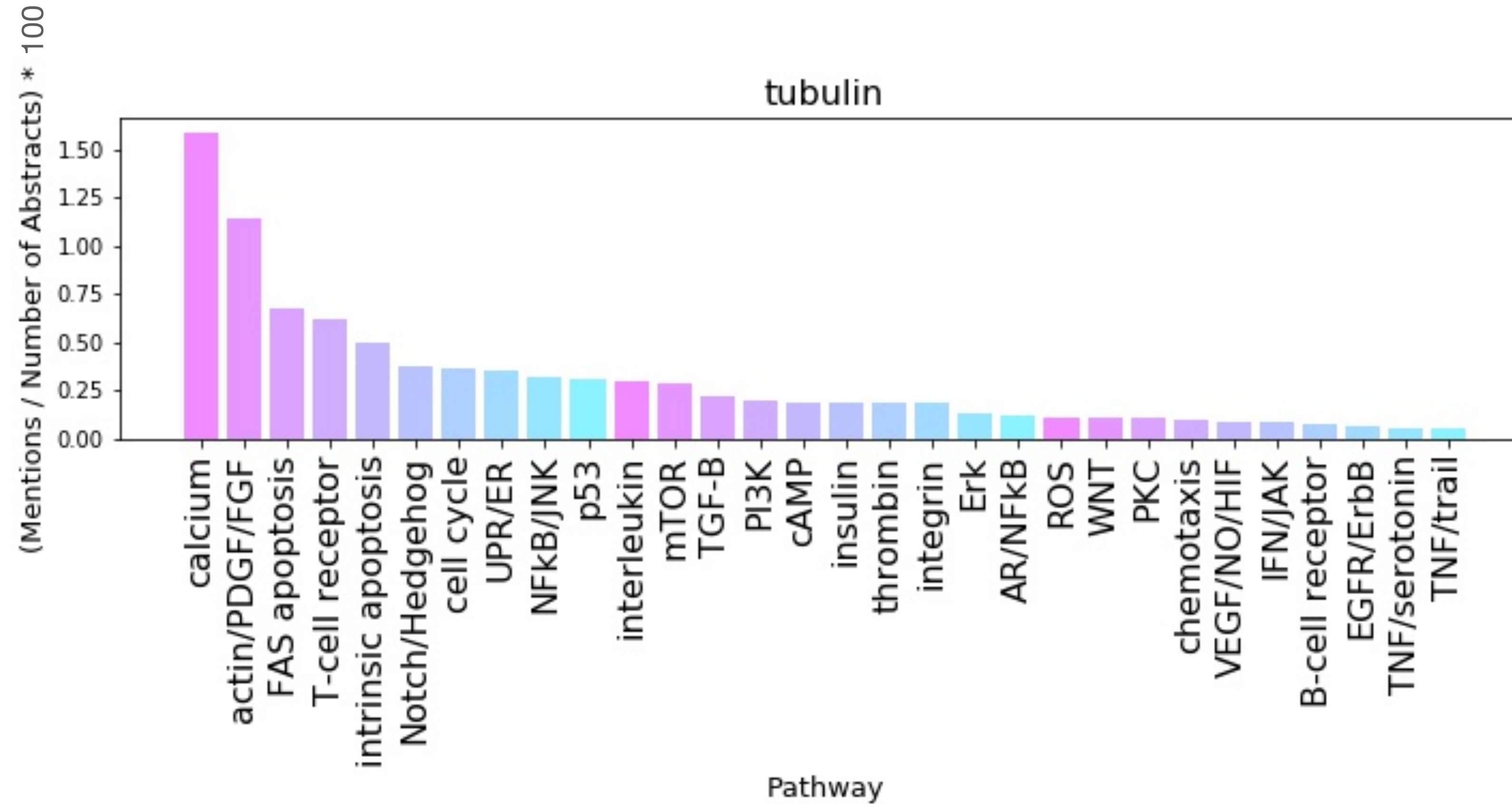
**Some clusters of
well-known related
pathways always
appear adjacent in
tSNE**



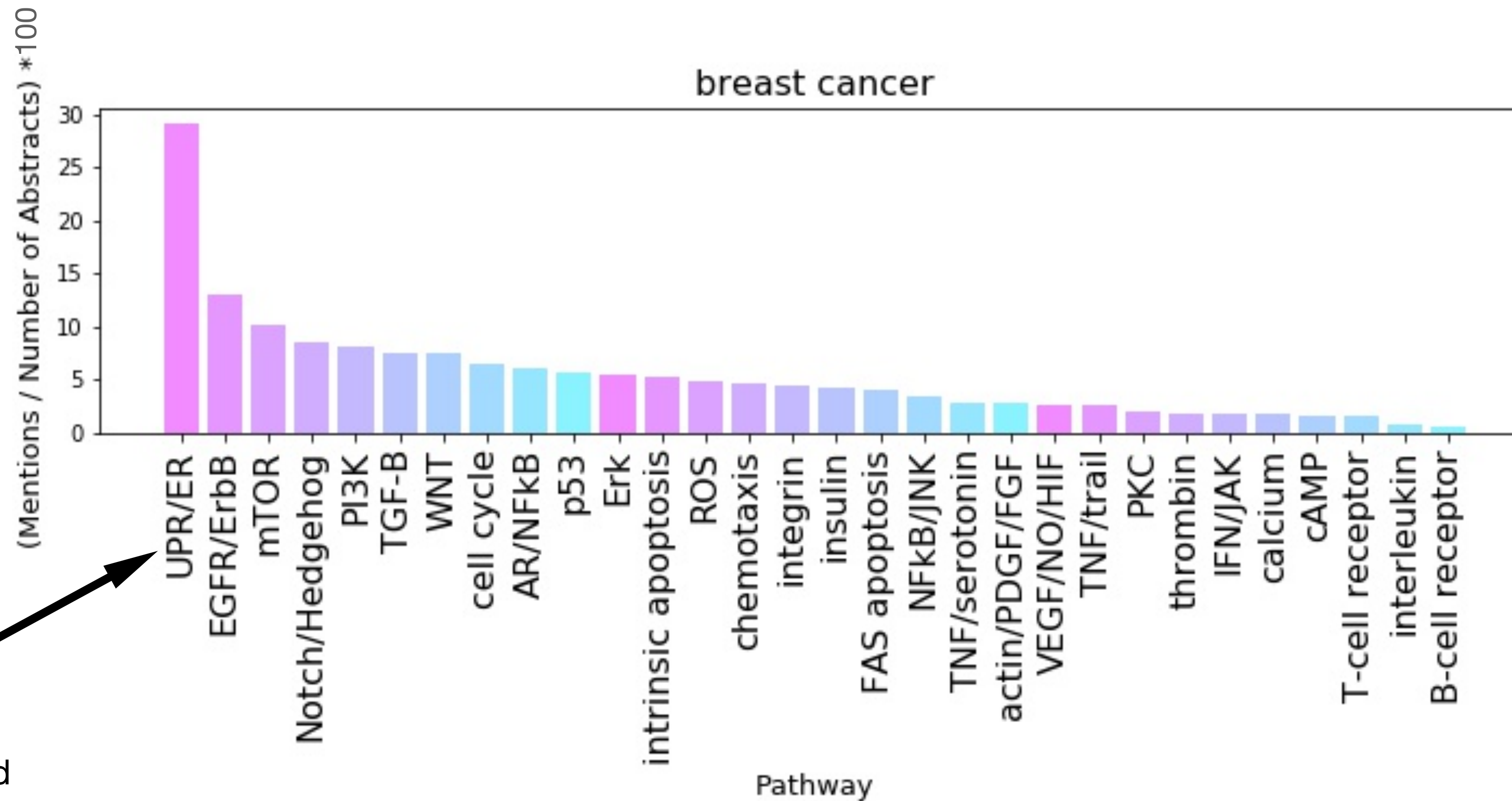
**Other clusters
represented
pathways that are
not as ‘canonical’ as
the others**



Looking up representation of gene tokens by cluster



Looking up representation of other tokens by cluster



Estrogen Receptor and
Endoplasmic Reticulum
have same acronym

Part II: Word2Vec Embedding

Corpus

>300k abstracts on human 'Cell Physiological Processes' and 'signaling' or 'pathway'

Word Embedding

List of tokens per sentence

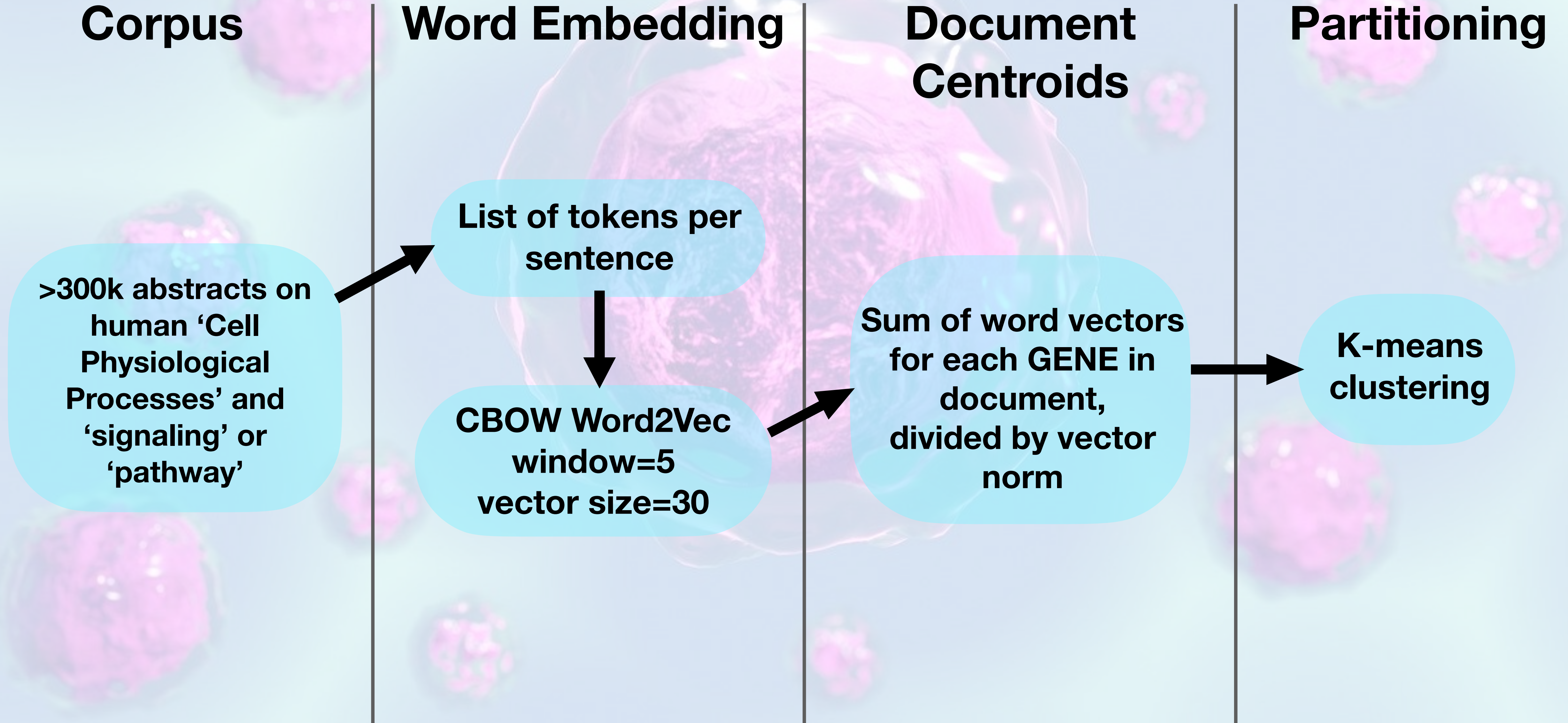
CBOW Word2Vec
window=5
vector size=30

Document Centroids

Sum of word vectors for each GENE in document, divided by vector norm

Partitioning

K-means clustering



Biology Analogies from Word2Vec

Activated by

Smad : Tgf-B :: STAT : JAK

Activates

cAMP : PKA :: calcium : CaMKII

Inhibited by

Nrf-2 : Keap-1 :: NF-kB : IKKa (2nd)

Interacts with

B-catenin : TCF :: Fos : JunB

Unit of

actin : F-actin :: tubulin : microtubule (2nd)

Fluid part

heart : blood :: brain : csf (2nd)

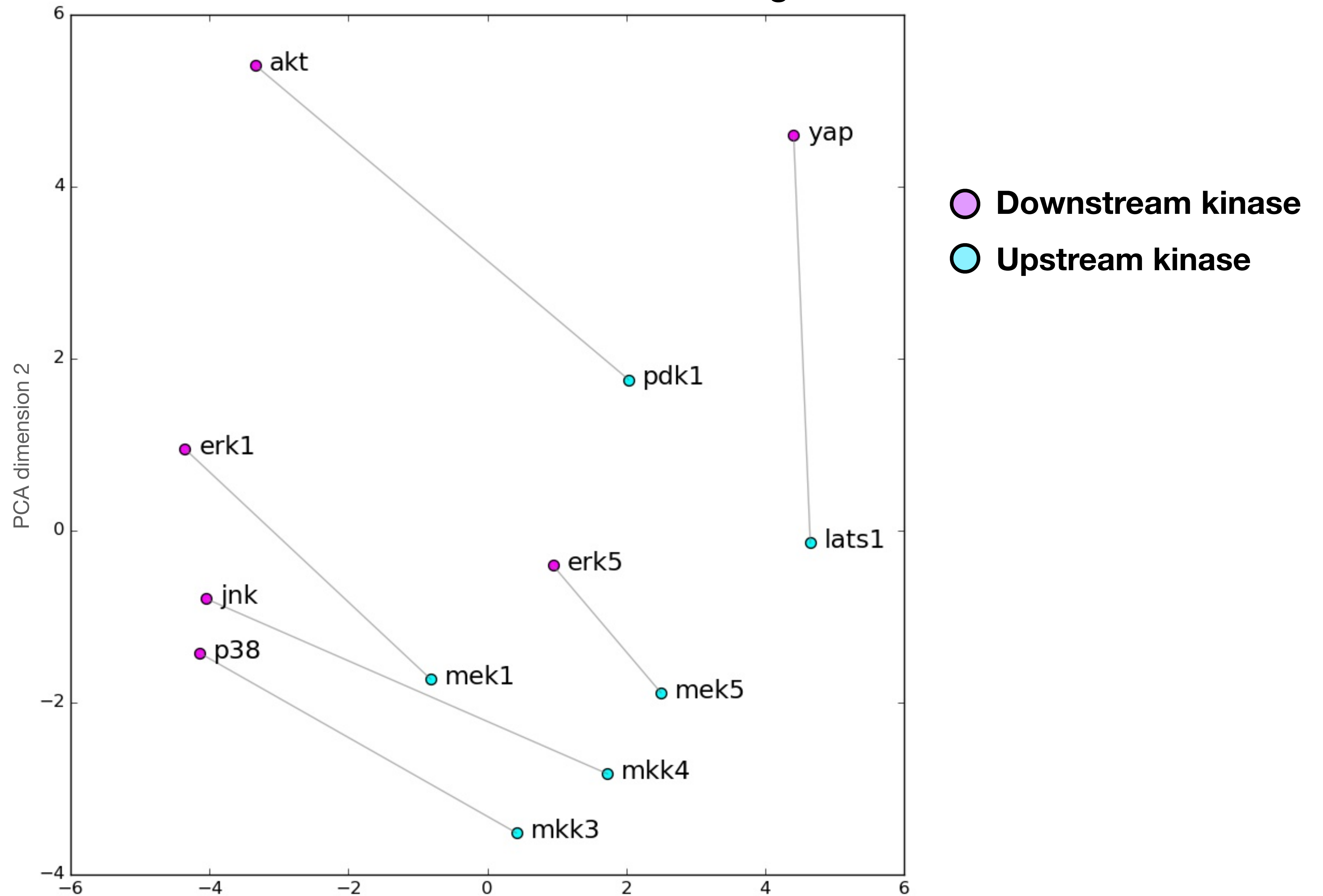
Hormone produced by

insulin : pancreas :: erythropoietin : fetal liver

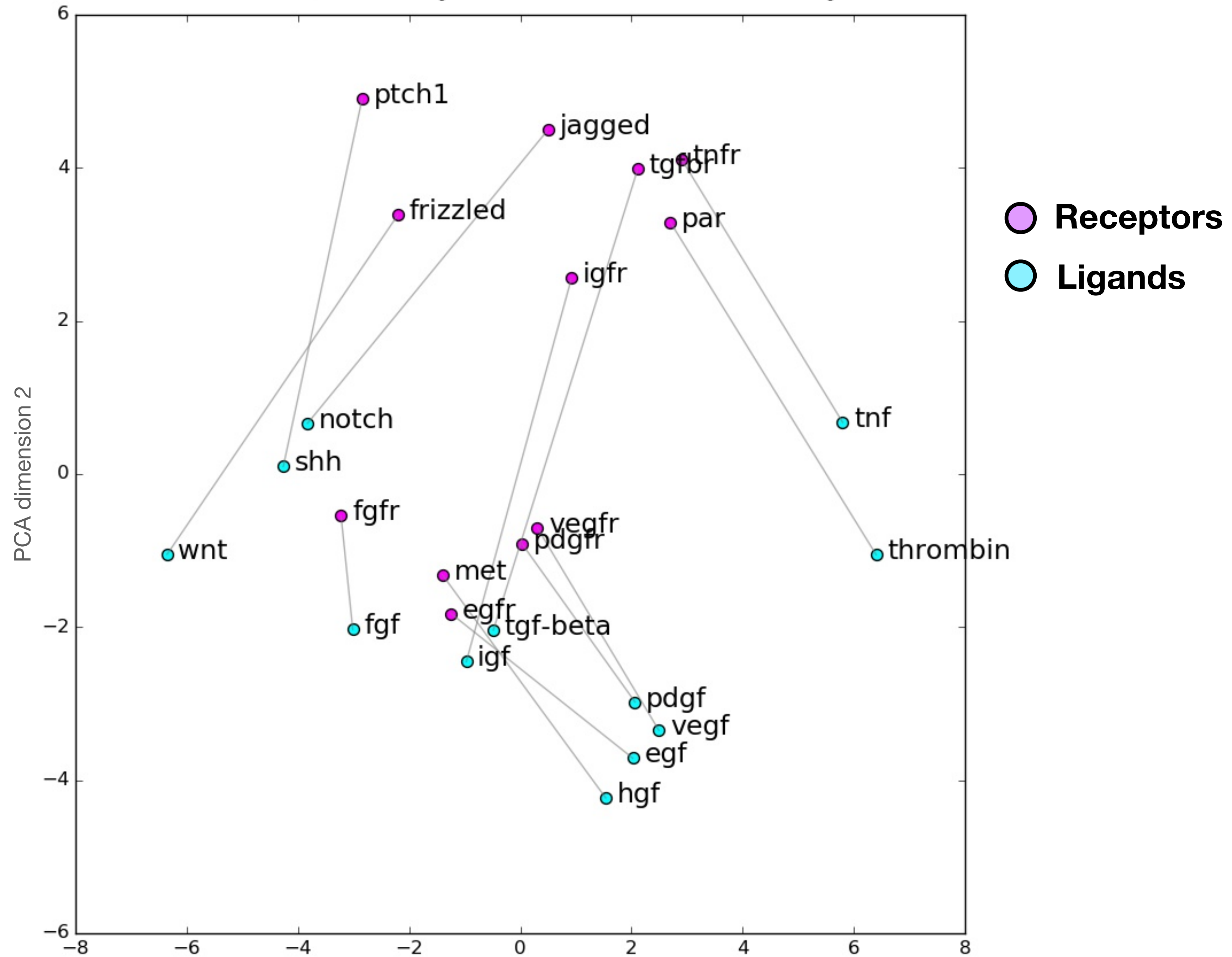
Binds to receptor

glucose : GLUT4 :: glutamate : AMPAR

Kinase Cascade Word Embeddings

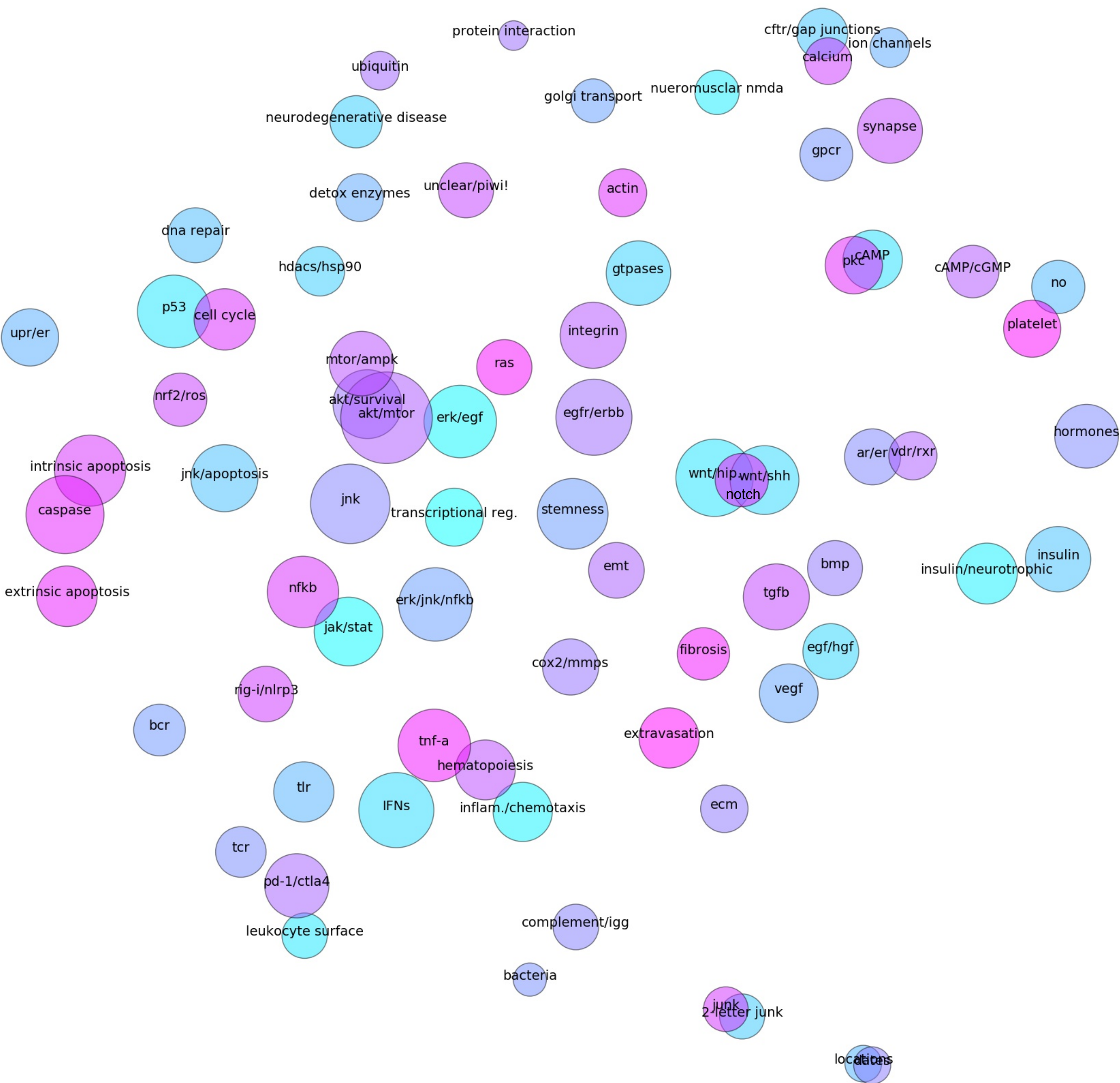


Receptor-Ligand Word Embeddings

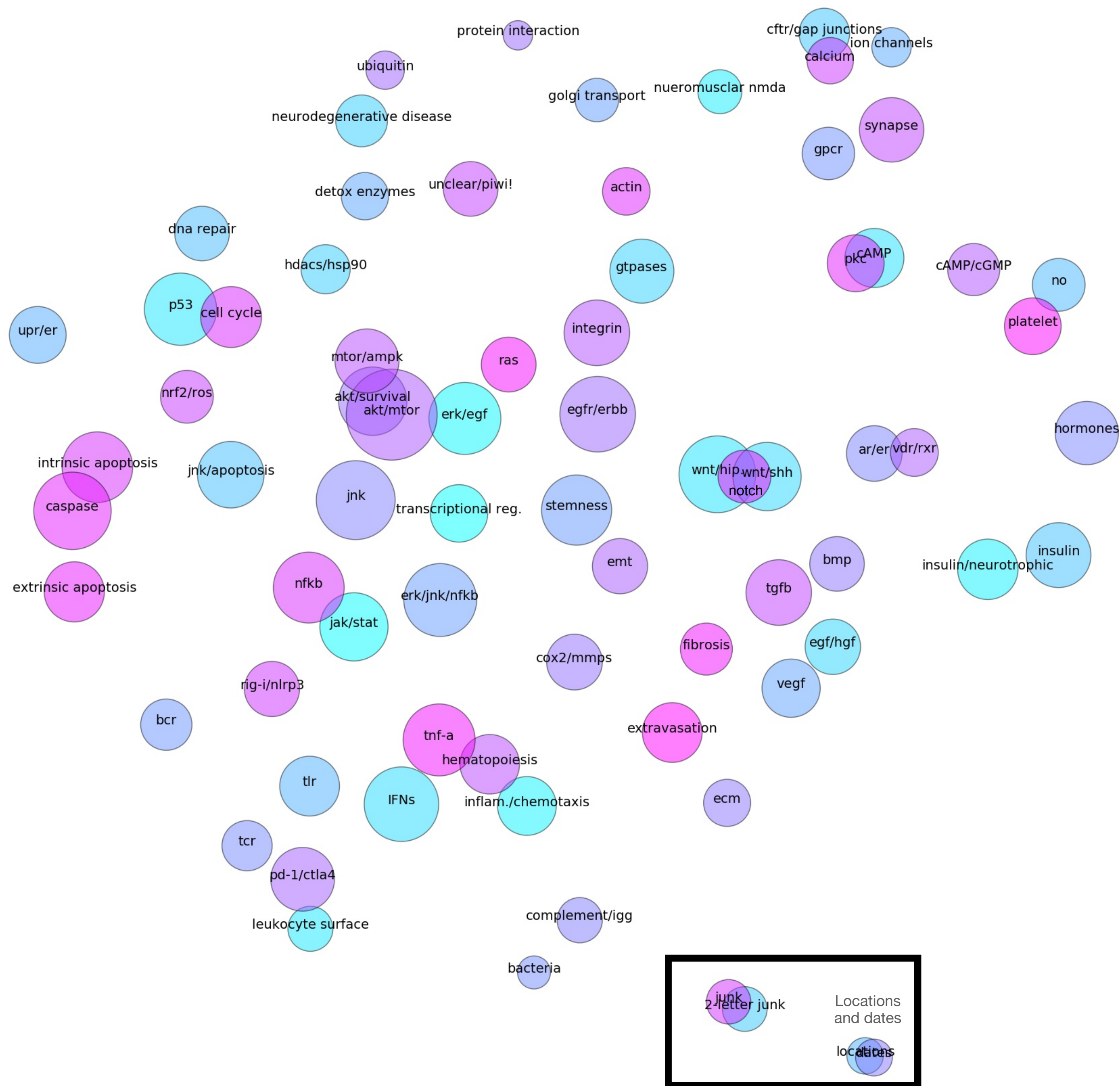


**More nuanced
clusters from
document centroids**

**Improvement over
NMF**



**Bad documents that
only contained
author info were
filtered out**



Future Directions

- App to search gene and token frequency for the Word2Vec-based clusters
- Entity linking tools to collapse gene and gene product aliases
- Sub-word embedding (learns partial word embeddings)
- Clustering techniques that allow greater variation in cluster size and shape
- Compare and upgrade to
 - BioConceptVec - CBOW word embeddings, entity linking, 30M abstracts
 - BioWordVec - Sub-word embeddings, 30M abstracts



Thank you!

