



# Biyomedikal Veri İçin Metin İşleme Tekniklerinin Kullanımı

Arzucan Özgür  
Boğaziçi Üniversitesi

Bilgisayar Mühendisliği Bölümü

[arzucan.ozgur@boun.edu.tr](mailto:arzucan.ozgur@boun.edu.tr)

Kuzeybatıda Yapay Öğrenme Yaz Okulu  
26.06.2019



# Biyomedikal Alanda Metin Tabanlı Veri

## ► Doğal Dildeki Metinler → Biyomedikal Metin Madenciliği

### Chitosan Oligosaccharide Exerts Anti-Allergic Effect against Shrimp Tropomyosin-Induced Food Allergy by Affecting Th1 and Th2 Cytokines.

Jiang T<sup>1,2</sup>, Ji H<sup>3</sup>, Zhang L<sup>4</sup>, Wang Y<sup>5</sup>, Zhou H<sup>6</sup>.

 Author information

#### Abstract

**BACKGROUND:** Shrimp-derived allergen has a serious impact on people's health. Chitosan oligosaccharide (COS) has anti-allergic action but its function on shrimp allergen-induced allergy and related molecular mechanisms remain unclear.

**METHODS:** COS and its degrees of polymerization (DP) were selected to interact with shrimp tropomyosin (TM) and IgE was measured. A mouse model of food allergy was established by receiving shrimp TM intraperitoneally. The models were treated with different concentrations of COS. Fecal and serum histamine, serum IgE, IgG1 and IgG2a, and inflammatory cytokines were measured.

**RESULTS:** The main products for COS were DP2-6 with the contents of 6, 40, 26, 16, and 4%, respectively, and reacted with shrimp TM increasingly when COS DP was increased. Severe symptoms of food allergy were observed in the TM group (diarrhea, anaphylactic response, and rectal temperature). In contrast, COS treatment improved these symptoms significantly ( $p < 0.05$ ). The sensitized mice were desensitized after they were treated with 1 mg/kg COS. COS treatment significantly reduced serum IgE and IgG1 levels, and increased IgG2a levels ( $p < 0.05$ ). COS consumption decreased fecal and serum histamine. COS treatment reduced Th2 cytokine (IL-4, IL-5, and IL-13) levels and increased the Th1 cytokine (IFN- $\gamma$ ) level ( $p < 0.05$ ).

**CONCLUSIONS:** COS showed anti-allergy properties by regulating the levels of Th1 and Th2 cytokines.

## ► DNA Dizisi → Genomik Veri İçin Metin Madenciliği

TTCAGGTGCATAAGACCTTGAC...

## ► Protein Dizisi ve İlaç Formülleri → İlaç Keşfi İçin Metin Madenciliği

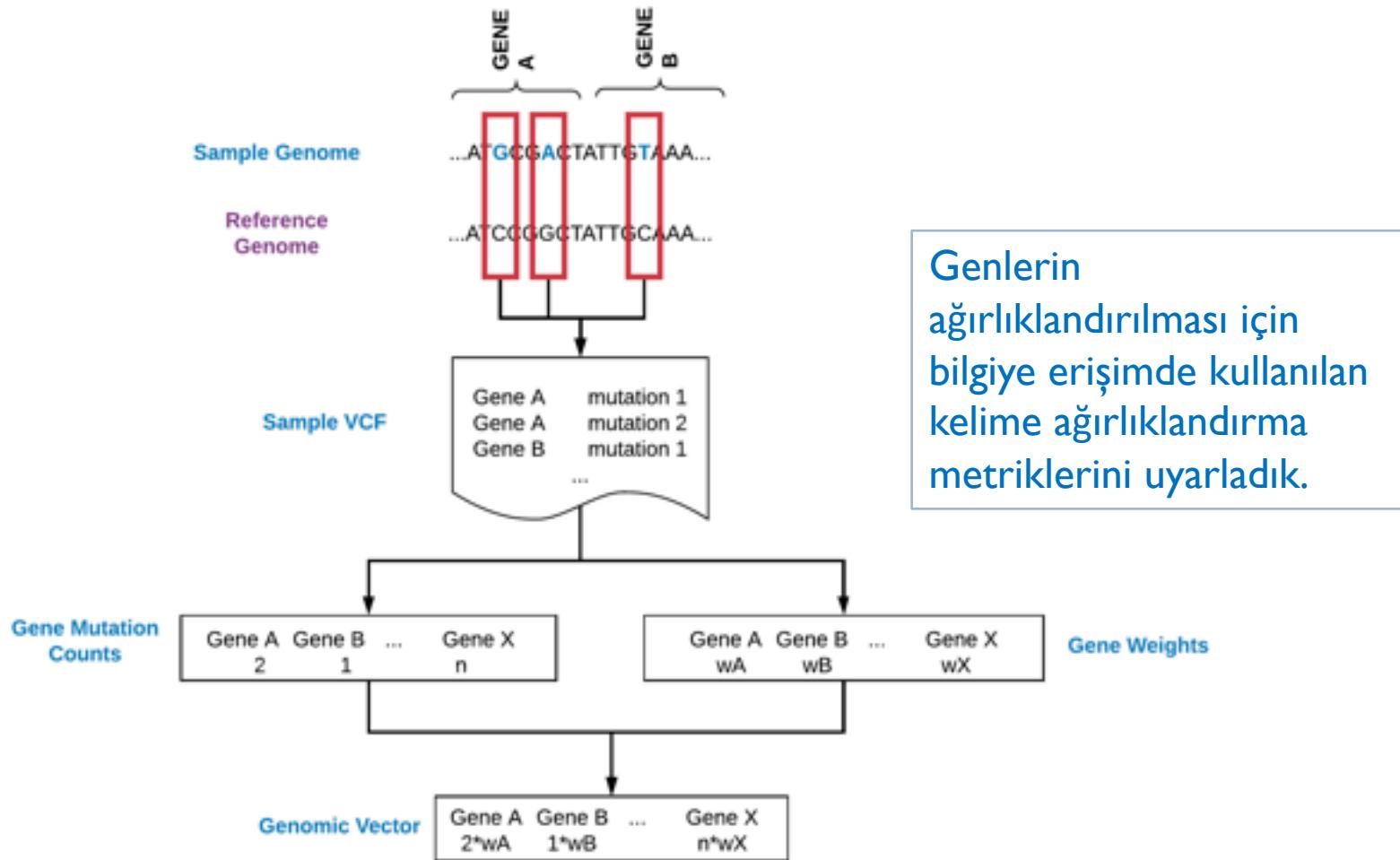
MELPNIMHPVAKLSTALAAALML...

CC1(C(N2C(SI)C(C2=O)NC(=O)...

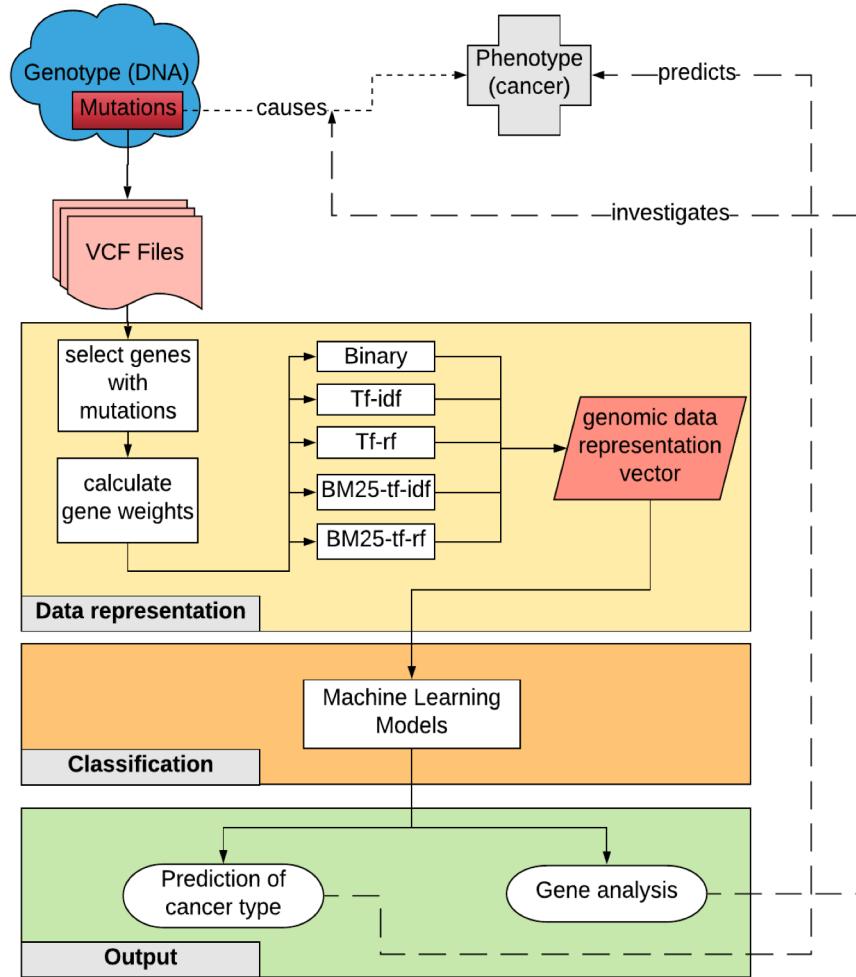
# Genomik Veri İçin Metin Madenciliği

İşbirliği: N. Özlem Özcan Şimşek ve Fikret Gürgen

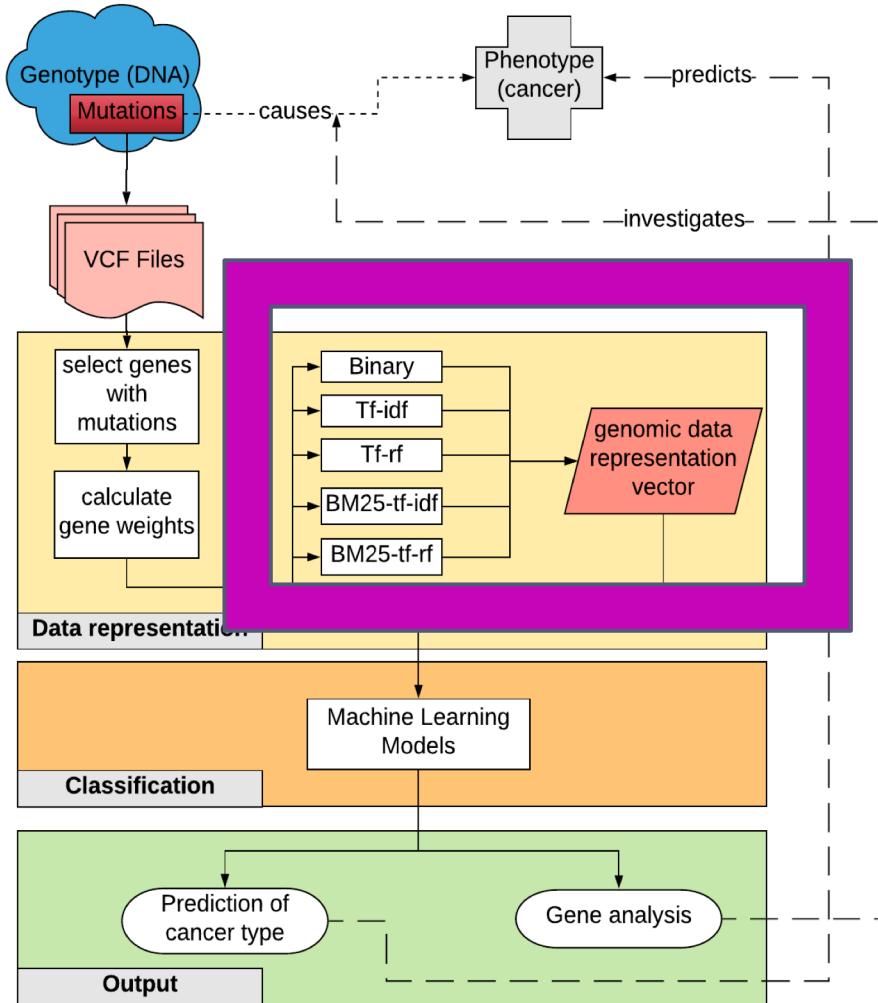
# DNA'daki Mutasyonlara Dayanan Hastalık Tahmini



# Sistemin Genel İşleyışı



# Veri Gösterim Modelleri



# İkili Gösterim

---

Bir gende mutasyon varsa 1 ile, yoksa 0 ile ifade ediliyor.

GEN-A (mutasyon var)	GEN-B (mutasyon yok)
1	0

## Kelime Frekansı – Ters Doküman Frekansı (TF-IDF)

- Bir gende ne kadar çok mutasyon varsa, o gen o kadar önemlidir (tf)
- Nadir mutasyonlar daha ayırdedididir (idf).

GEN-A (genel mut)	GEN-B (nadir mut)
0.0056	0.035

$$tf\text{-}idf_{g,s} = tf_{g,s} * idf_g$$

mutasyon sayısı

$\log(N/G)$  (**N**: hasta sayısı; **G**: g geninde mutasyon olan hasta sayısı)

## Kelime Frekansı – İlgililik Frekansı (TF-RF)

- Sınıf bilgisi kullanan denetimli bir yaklaşım.
- Bir gende belirli bir hastalık türünde diğer hastalık türlerine göre daha fazla mutasyon varsa, o gen daha ayırdedidir.

	GEN-A	GEN-B
Hastalık-X	0.0056	0.035
Hastalık-Y	1.0056	0.005

$$tf\text{-}rf_{g,s} = tf_{g,s} * rf_{g,c}$$

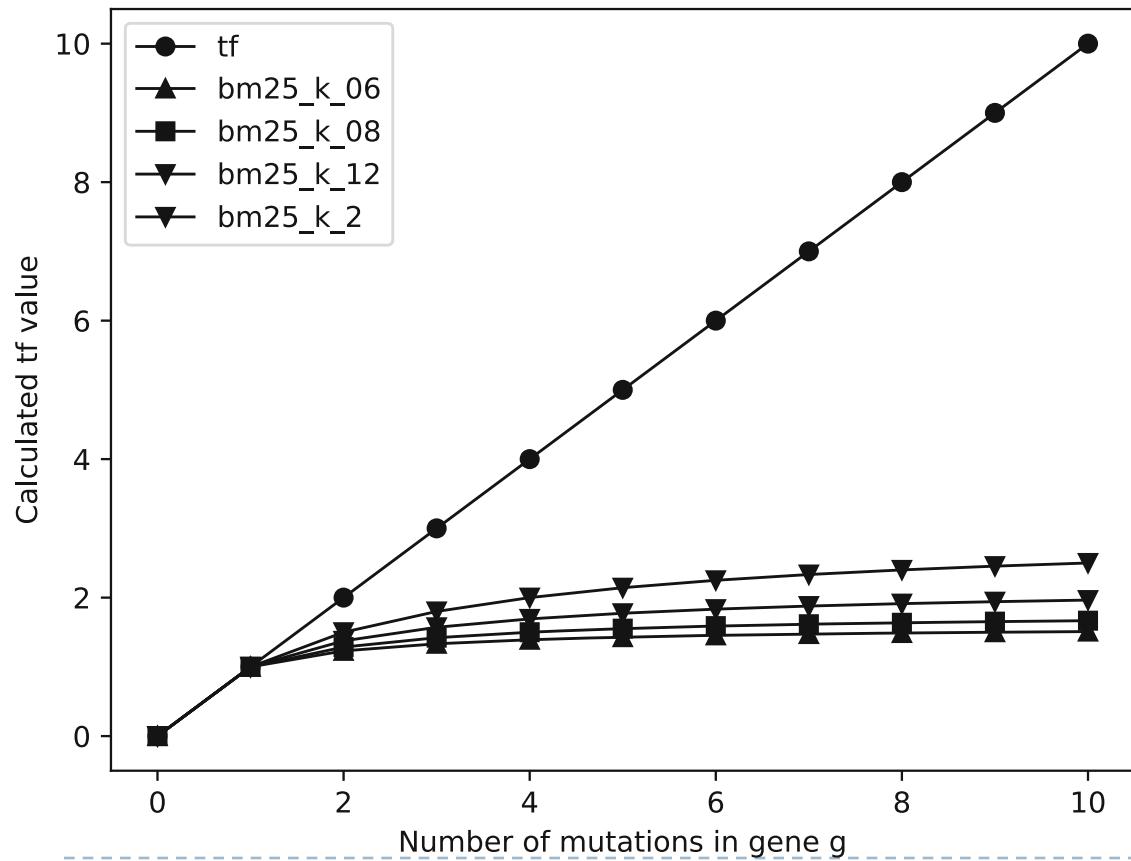
mutasyon sayısı

$$\log(2+a/\max(l,b))$$

a: c sınıfında g geninde mutasyon olan hasta sayısı  
b: diğer sınıflarda g geninde mutasyon olan hasta sayısı

# BM25

tf → BM25-tf



Yumuşatma  
(sınırlama) etkisi

# BM25-TF-IDF & BM25-TF-RF

$$BM25-tf-idf_{g,s} = BM25-tf_{g,s} * idf_g$$

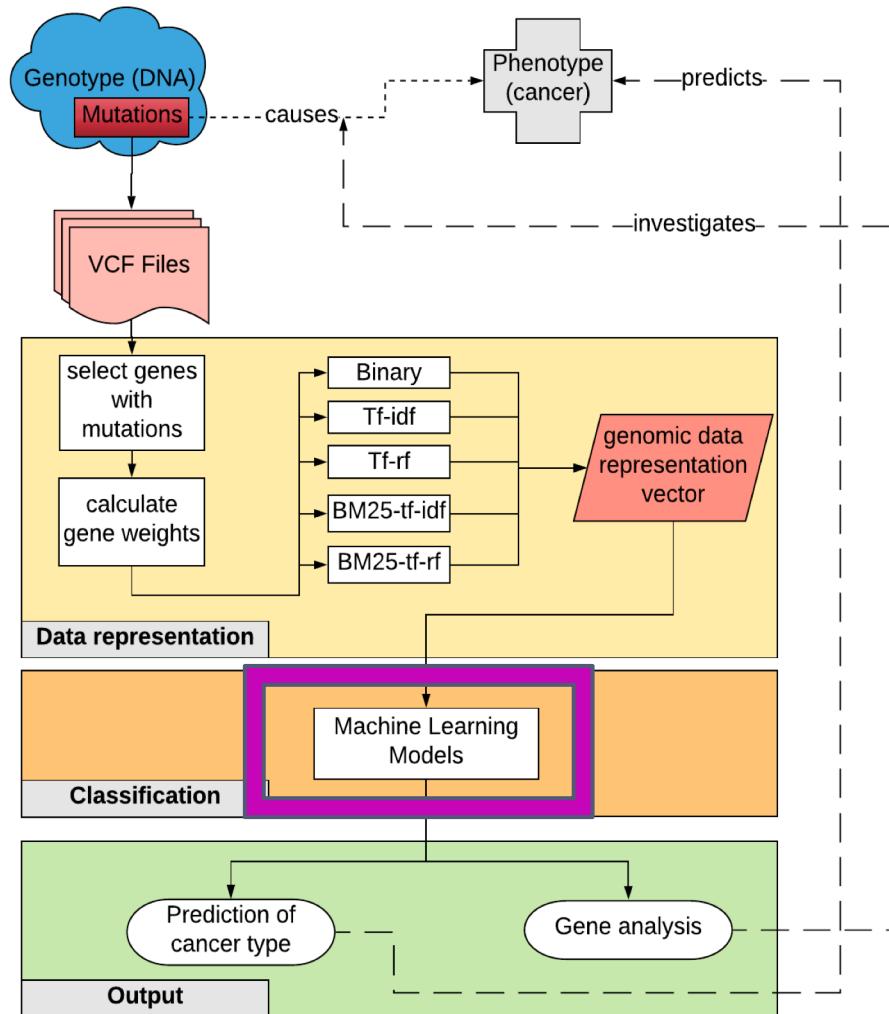
## yumuşatılmış mutasyon sayısı

# gen ağırlığı

$$BM25-tf \cdot rf_{g,s} = BM25-tf_{g,s} * rf_{g,c}$$

yumuşatılmış  
mutasyon sayısı denetimli gen  
ağırlığı

# Sınıflandırma



# Sınıflandırma

---

- Naive Bayes (NB)
- K En Yakın Komşu (KNN)
- Destek Vektör Makineleri (SVM)
- Yapısal Bağıntı (Logistic Regression - LR)
- Algılayıcı (Perceptron)
- Çok Katmanlı Yapay Sinir Ağları (NN)

# Veri Kümesi

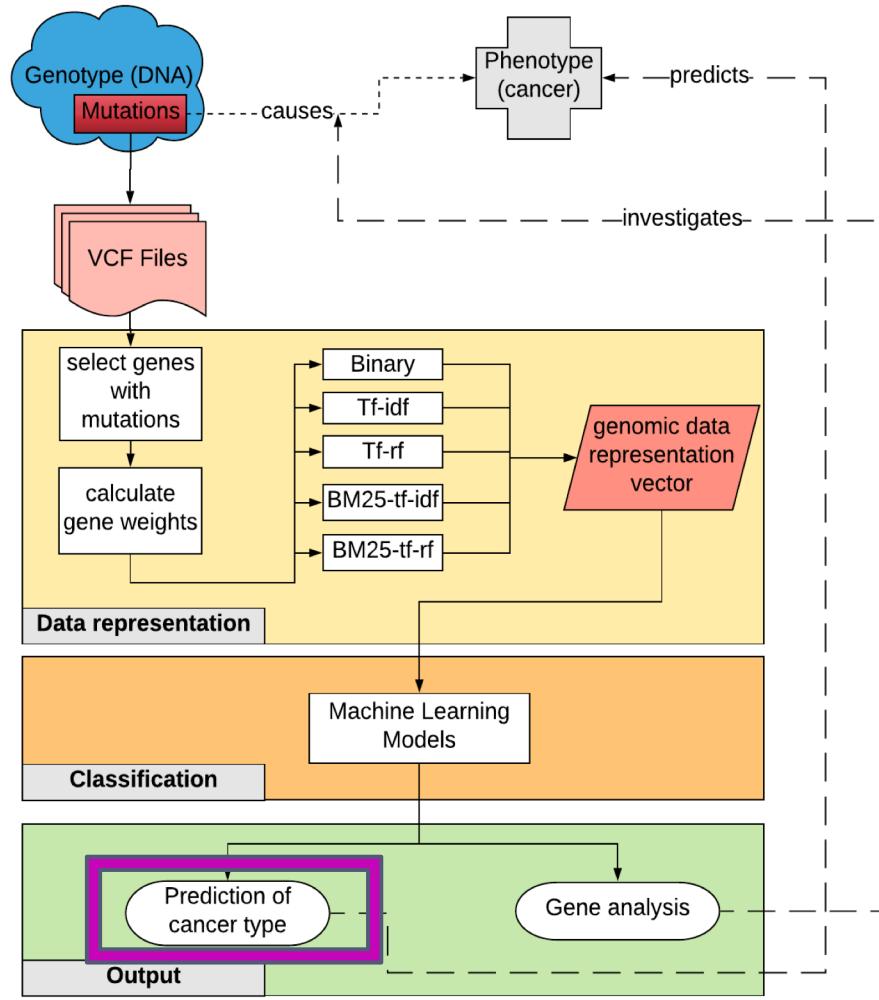
## TCGA (The Cancer Gene Atlas)

Cancer type	Sample count
Lung	1232
Breast	1080
Brain	1028
Kidney	734
Colorectal	656
Thyroid	504
Prostate	503
Skin	472
Stomach	441
Liver	378

Training & Validation  
(80%)

Test  
(20%)

# Kanser Türü Tahmini



# Sonuçlar

Bilinen ve tüm genlerin karşılaştırılması

Gen Kümesi	Doğruluk	F-ölçütü	Yanlış Pozitif Oranı
Bilinen	$36.74 \pm 0.56$	$36.62 \pm 0.83$	$10.01 \pm 0.10$
Tüm	$68.50 \pm 0.48$	$69.01 \pm 0.01$	$4.07 \pm 0.09$

\* LR sonuçları

# Sonuçlar

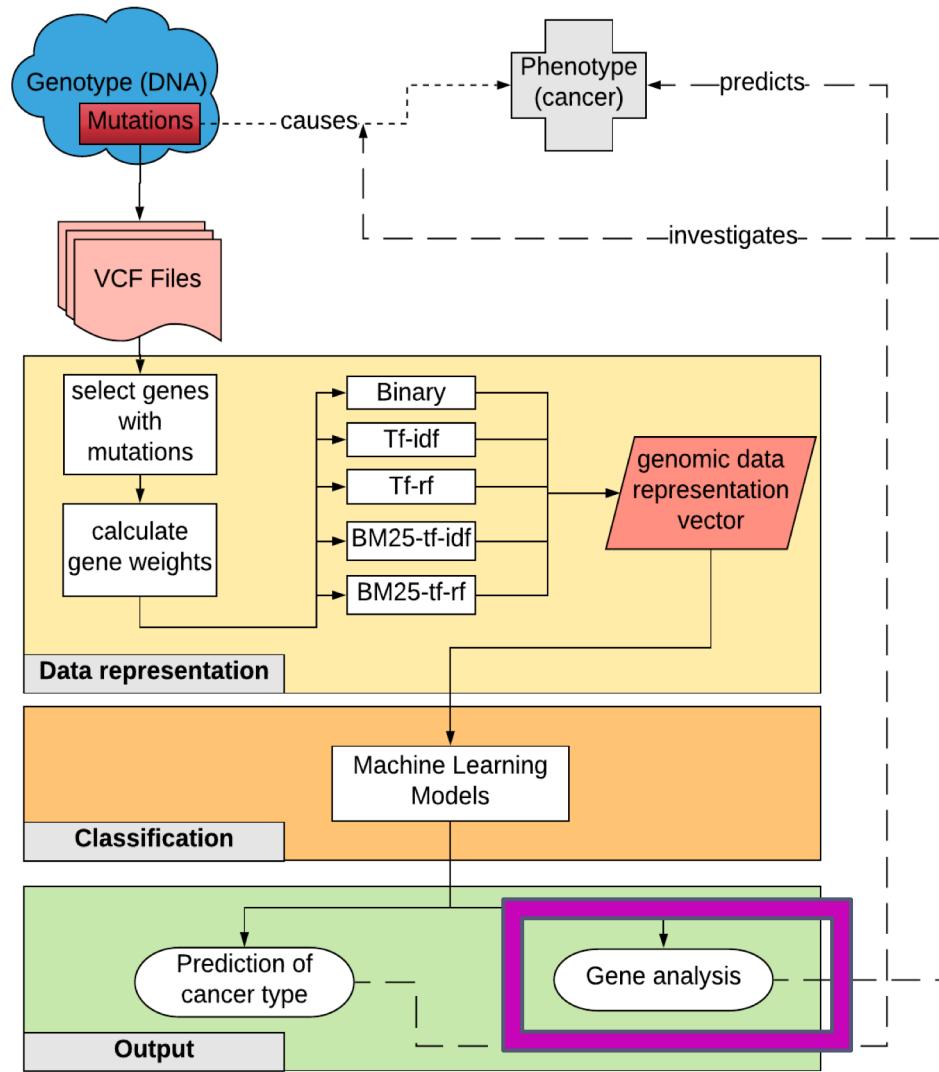
Veri gösterim modellerinin karşılaştırılması

Veri gösterimi	Doğruluk	F-ölçütü	Yanlış Pozitif Oranı
İkili	$69.00 \pm 0.76$	$69.52 \pm 0.70$	$3.65 \pm 0.17$
Tf-idf	$62.91 \pm 0.79$	$63.32 \pm 0.70$	$4.00 \pm 0.10$
Tf-rf	$74.13 \pm 1.33$	$74.17 \pm 1.47$	$3.07 \pm 0.24$
BM25-tf-idf	$68.18 \pm 1.83$	$68.79 \pm 1.28$	$4.07 \pm 0.54$
BM25-tf-rf	$76.44 \pm 0.66$	$76.95 \pm 0.68$	$2.75 \pm 0.13$
C-score	$73.74 \pm 0.88$	$74.07 \pm 0.73$	$3.27 \pm 0.24$

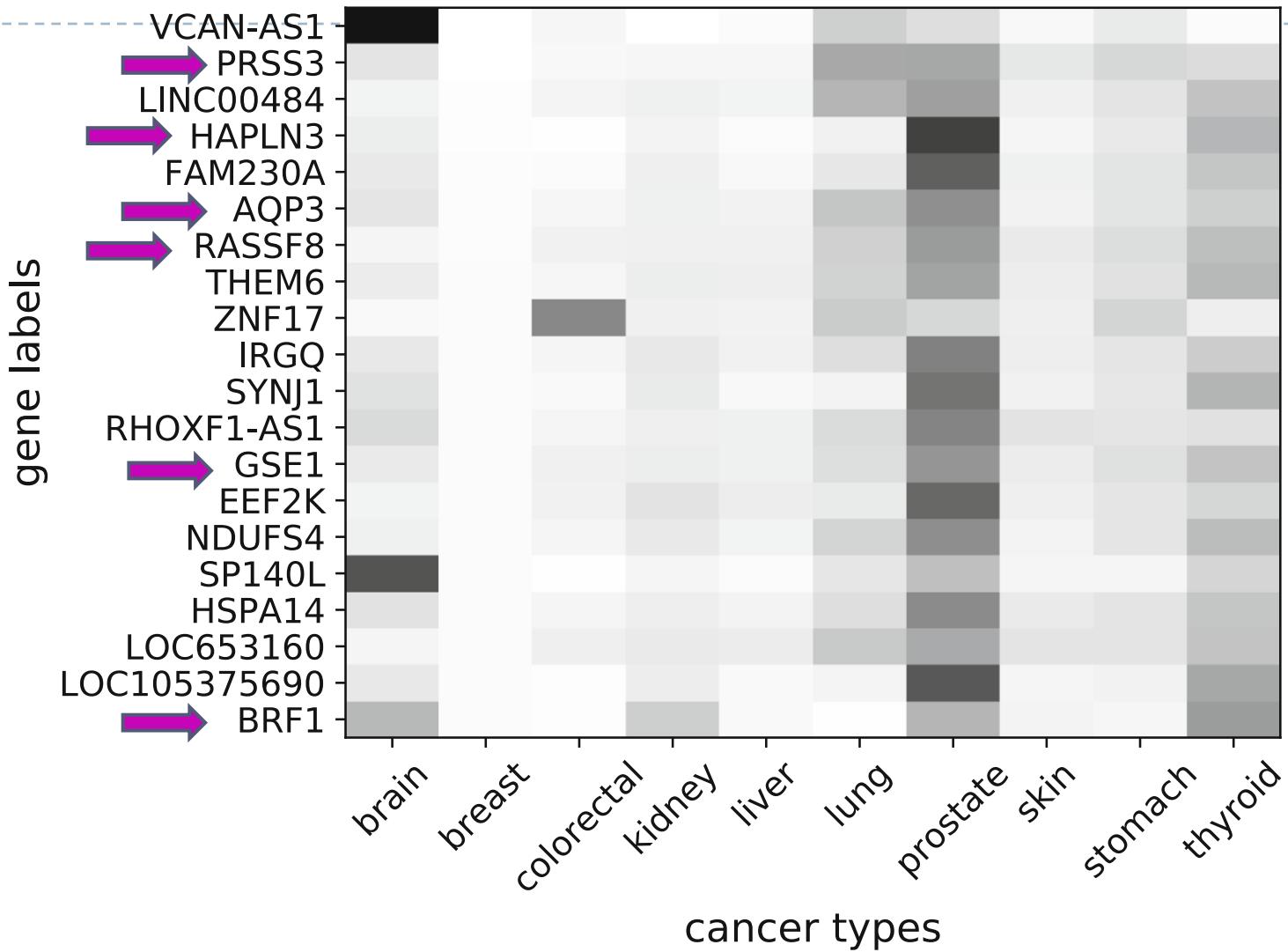
\* NN sonuçları

Statistical significance 95% conf. interval p-value: 0.0001

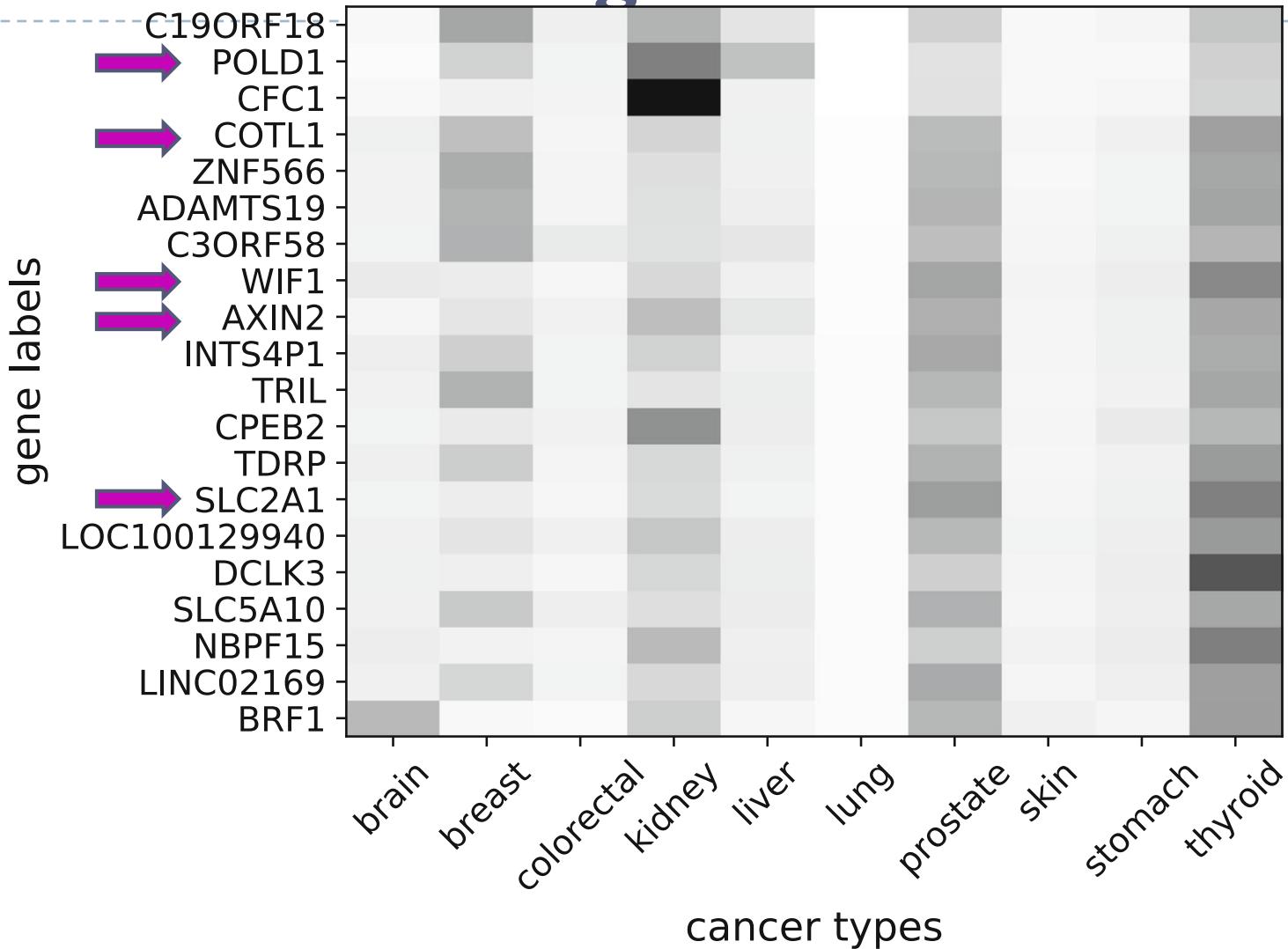
# Gen Analizi



# Gen Analizi – Meme kanseri



# Gen Analizi – Akciğer Kanseri

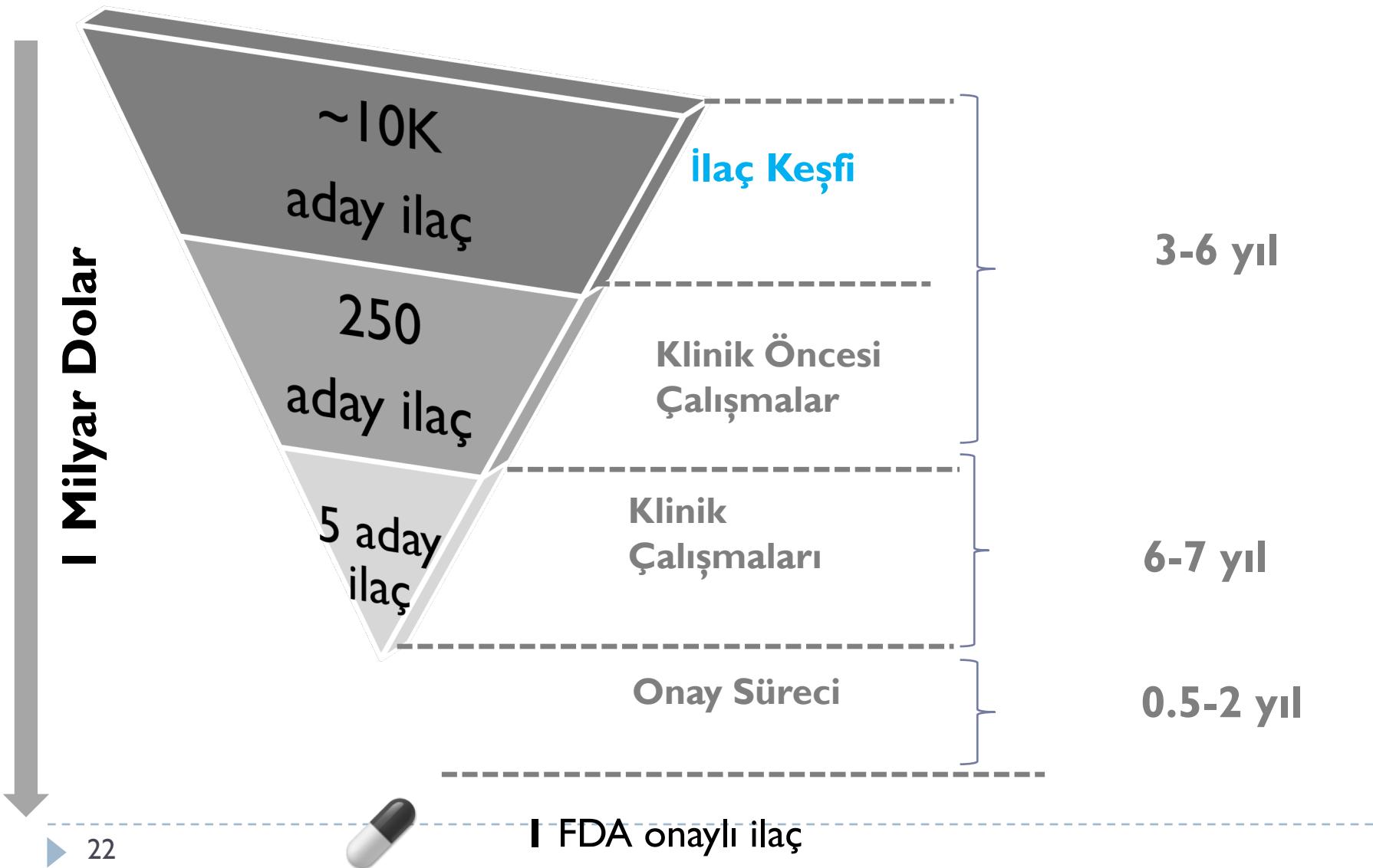


# İlaç Keşfi İçin Metin Madenciliği

İşbirliği: Hakime Öztürk ve Elif Özkırımlı

# Motivasyon:

İlaç geliştirme uzun, pahalı ve zor bir süreç



# Bu süreci nasıl kolaylaştırabiliriz

İlaç adaylarını önceliklendirerek arama uzayını daraltabiliriz.



PubChem > DrugBank



UniProt > DrugBank

Mevcut ilaçlar için yeni hedefler belirleyebiliriz.

2



Yeni tedavi yöntemleri

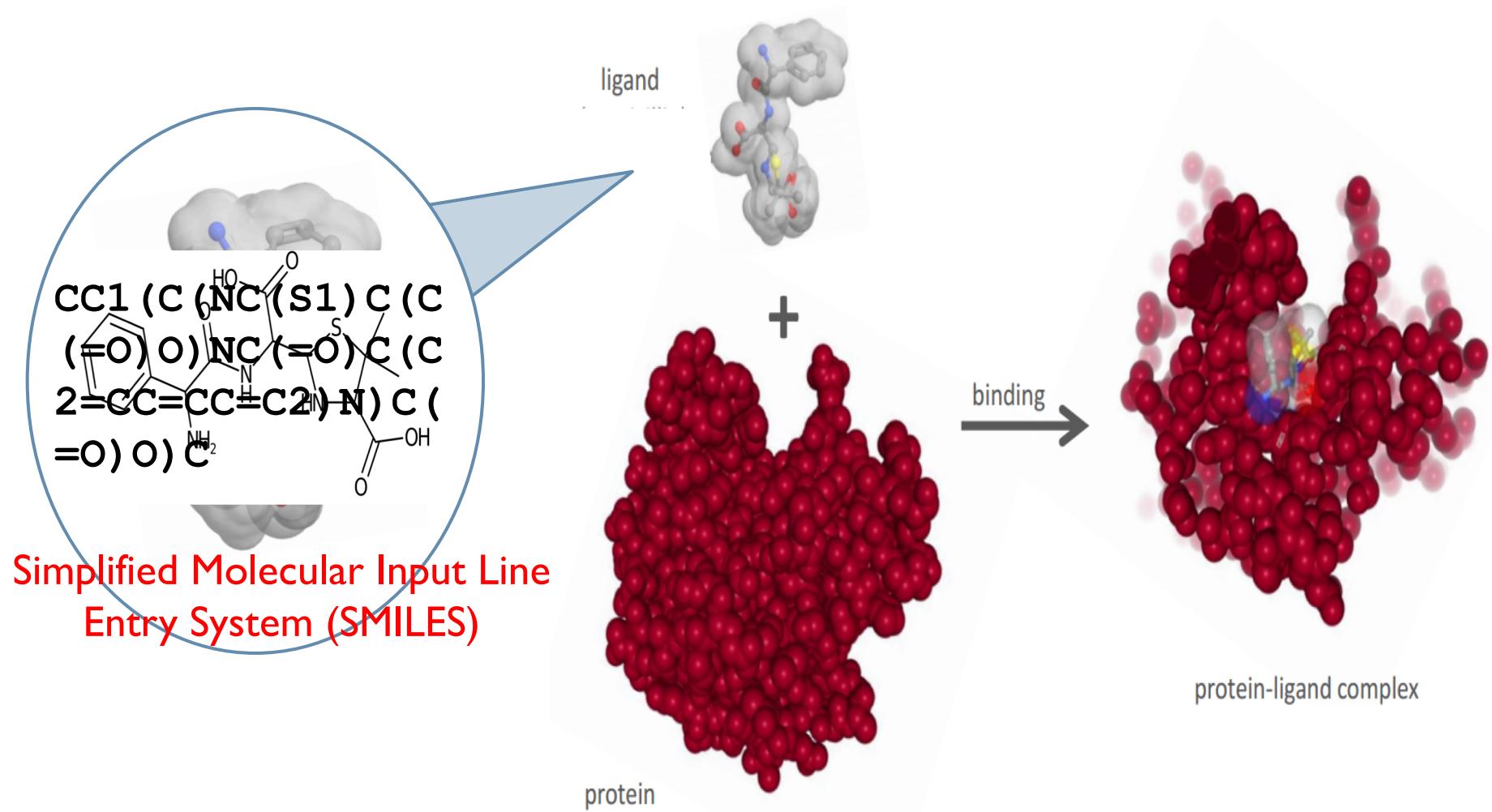
Bilinen ilaçlar

BindingDB  
ChEMBL

## Yapay öğrenme teknikleri kullanarak



# Protein ve ligandların metin tabanlı gösterimi



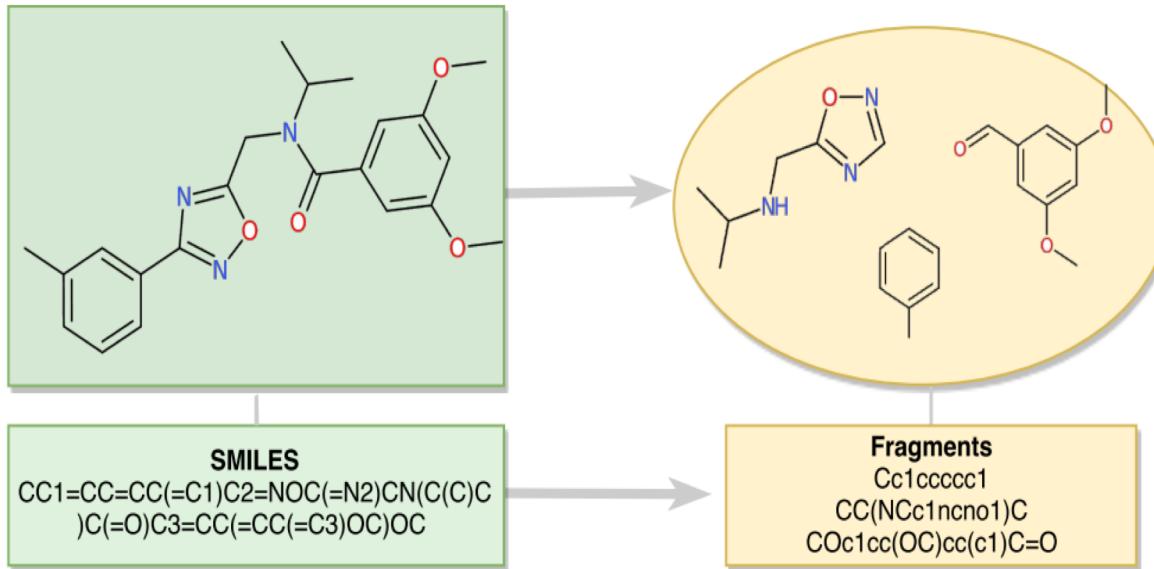
# SMILES metin tabanlı bir gösterim

---

- Metin işleme yöntemleri kullanabiliriz
- Ligand gösterimi için iki yaklaşım kullandık
  - TF-IDF tabanlı
  - Dağıtık gösterim tabanlı



## ▶ Varsayımlı: Kelimelerden oluşan doküman



- Kelimelerin ne olduğunu bilmiyoruz

# Doküman olarak **SMILES**

---



**SMILES:**

COCl=C(C=CC(=Cl)C=O)O

# Doküman olarak **SMILES**

---

**SMILES:**

8-karakterlik **LINGO**



**COCl=C(C=CC(=Cl)C=O)O**

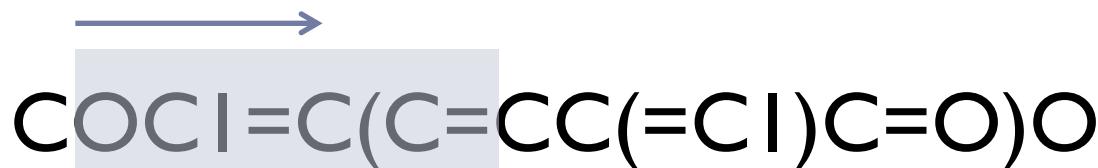
**Kimyasal Kelimeler:** COCl=C(C

# Doküman olarak **SMILES**

---

**SMILES:**

8-karakterlik **LINGO**



**Kimyasal Kelimeler:** COCl=C(C=OCl)C=

# Doküman olarak **SMILES**

---

**SMILES:**

8-karakterlik **LINGO**



**Kimyasal Kelimeler:** COCl=C(C  
OCl=C(C=  
Cl=C(C=C

# Doküman olarak **SMILES**

---

## **SMILES:**

8-karakterlik **LINGO**



**Kimyasal Kelimeler:** COClI=C(C  
OCI=C(C=  
CI=C(C=C  
I=C(C=CC

# Doküman olarak **SMILES**

---

## **SMILES:**

8-karakterlik **LINGO**

COCl=C(C=CC(=Cl)C=O)O

**Kimyasal Kelimeler:** COCl=C

OCl=C(C=

Cl=C(C=C

I=C(C=CC

...

Cl)C=O)O

# Kelime Frekansı – Ters Doküman Frekansı (TF-IDF)

---

- Kimyasal kelimeleri ağırlıklandırmak için kullandık
  - **TF**: Kimyasal kelimenin SMILES'da geçme frekansı
  - **IDF**: Kimyasal kelimenin derlemde geçtiği SMILES formülü sayısı ile ters orantılı (varsayımlı: nadir geçen kelimeler daha ayırdedici)



# Kelime frekansı (TF)

Kimyasal kelimenin bir SMILES formülünde geçme sıklığı:

Kimyasal kelimeler	ligand (SMILES)			
	SMILES <sub>1</sub>	SMILES <sub>2</sub>	SMILES <sub>3</sub>	SMILES <sub>4</sub>
=CC=CC=C	3	6	4	2
Cl)NC(=O		2		2
COCl=C(C	0		0	0
2=C(C(=N	0			0

# Ters Doküman Frekansı (IDF)

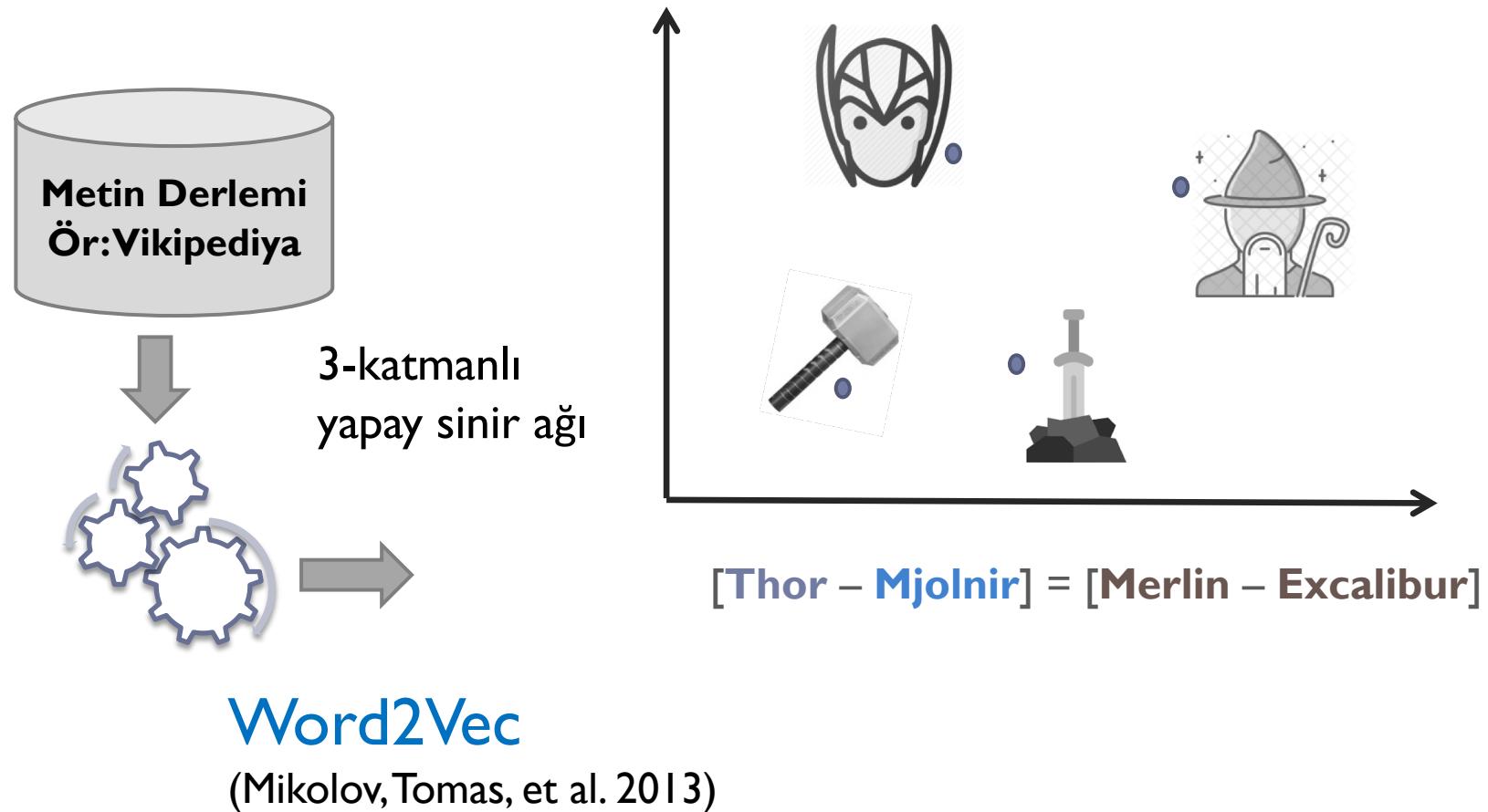
Chemical words	ligand (SMILE S)	SMILES <sub>1</sub>	SMILES <sub>2</sub>	SMILES <sub>3</sub>	SMILES <sub>4</sub>
=CC=CC=C		3	6	4	2
TÜM ligandlarda geçiyor.					
Cl)NC(=O			2		2
COCl=C(C		0		0	0
Tek bir ligandda geçiyor					
2=C(C(=N		0			0

Ligandları TF-IDF tabanlı vektörler olarak ifade ettik.

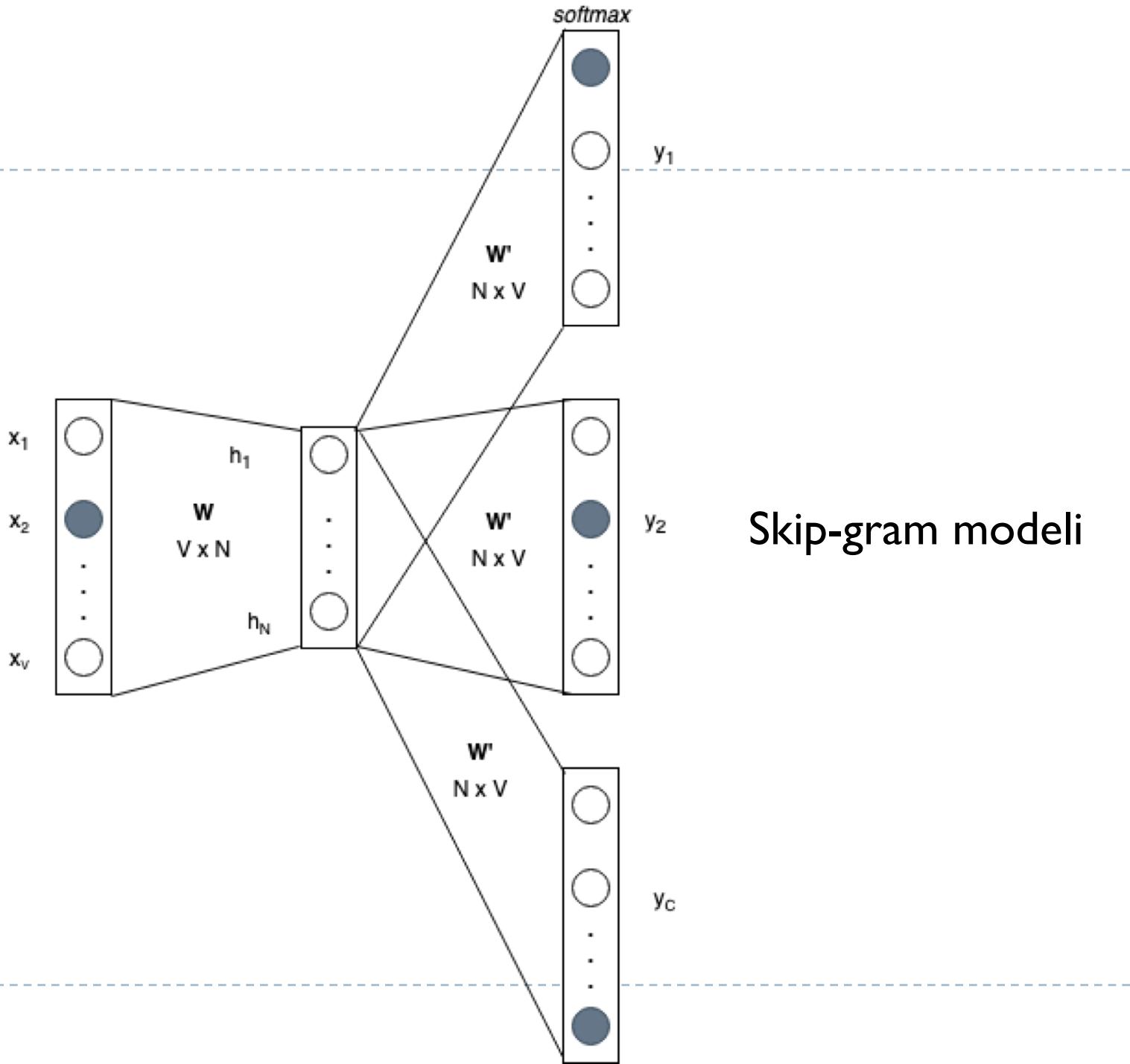
TF-IDF vektörleri arasındaki kosinüs benzerliğini k en yakın komşu tabanlı algoritma ile kullandık.  
İki boyutlu ligand gösterimi ile aynı başarıyı daha az hesaplama gücü ile elde ettik.

- Öztürk, H., Özkırımlı, E., Özgür, A. **A comparative study of SMILES-based compound similarity functions for drug-target interaction prediction.** *BMC Bioinformatics*, 17:128, 2016.

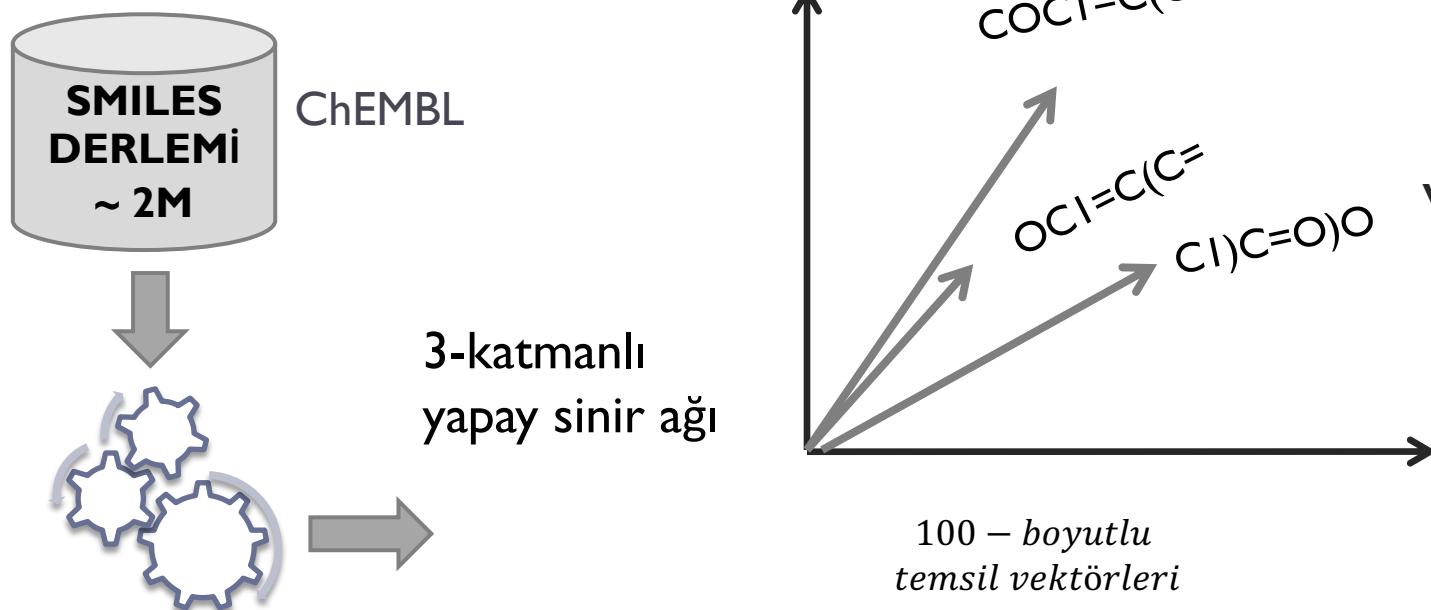
# Dağıtık Kelime Temsil Modeli



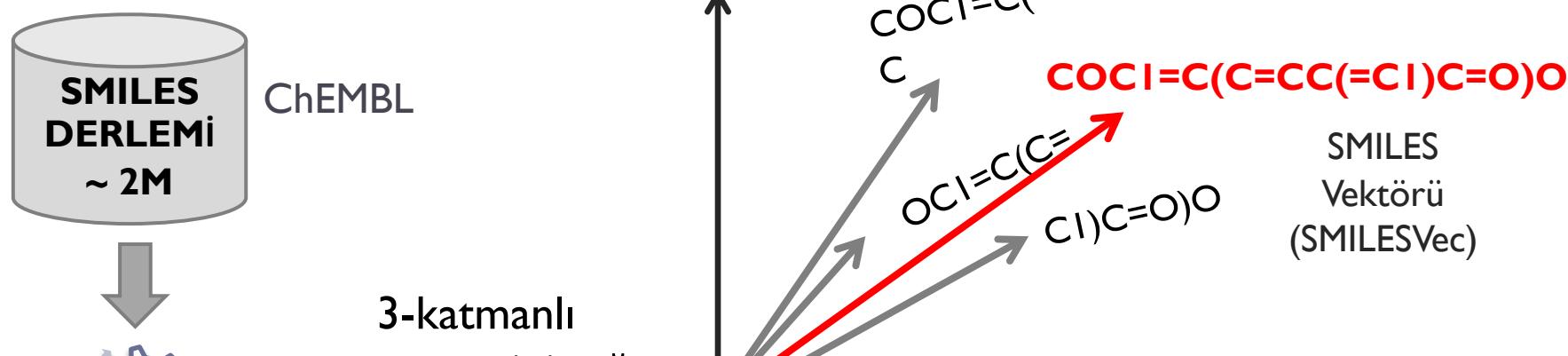
Word2Vec  
(Mikolov, Tomas, et al. 2013)



# SMILESVec: Dağıtık ligand temsili



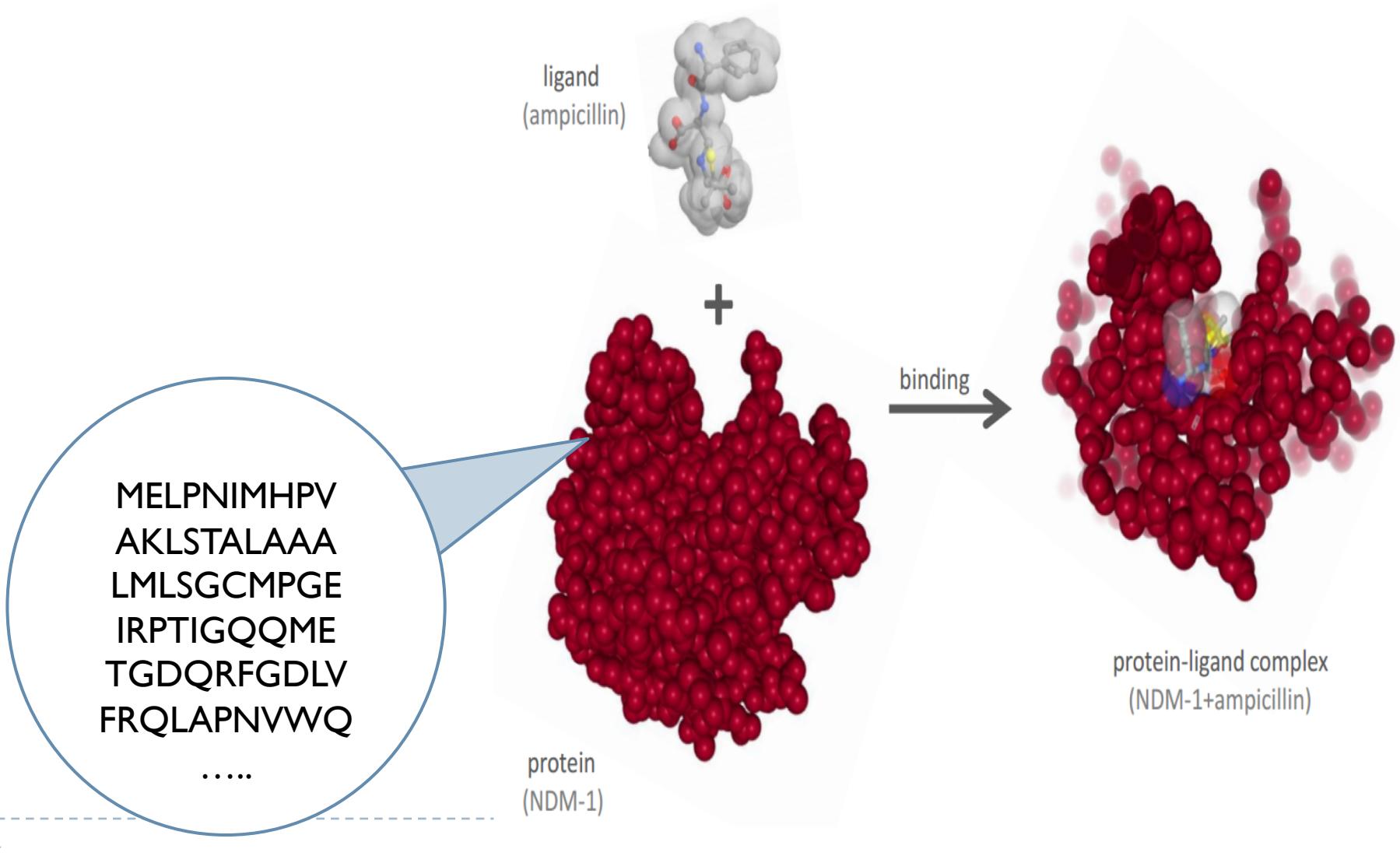
# SMILESVec: Dağıtık ligand temsili



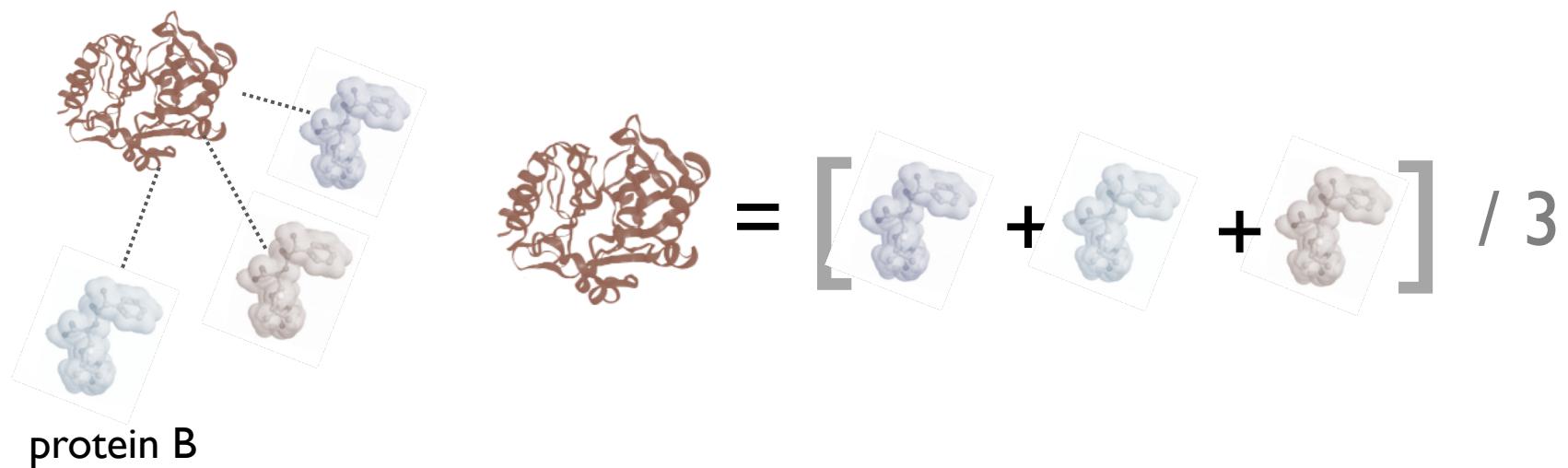
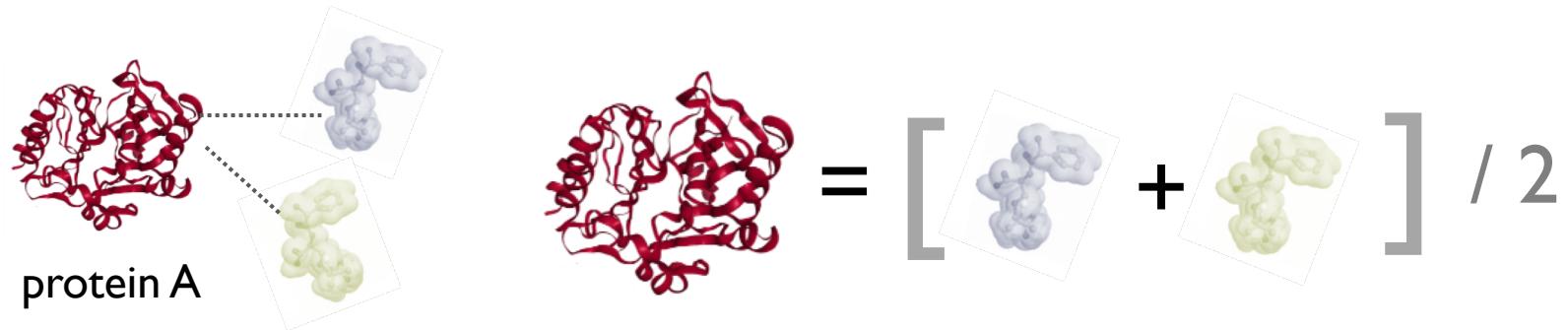
$$\text{SMILESVec} = \text{vector(ligand)} = \frac{\sum_{k=1}^n \text{vector(word}_k\text{)}}{n}$$

*n, kimyasal kelimelerin sayısı*

# Protein temsili



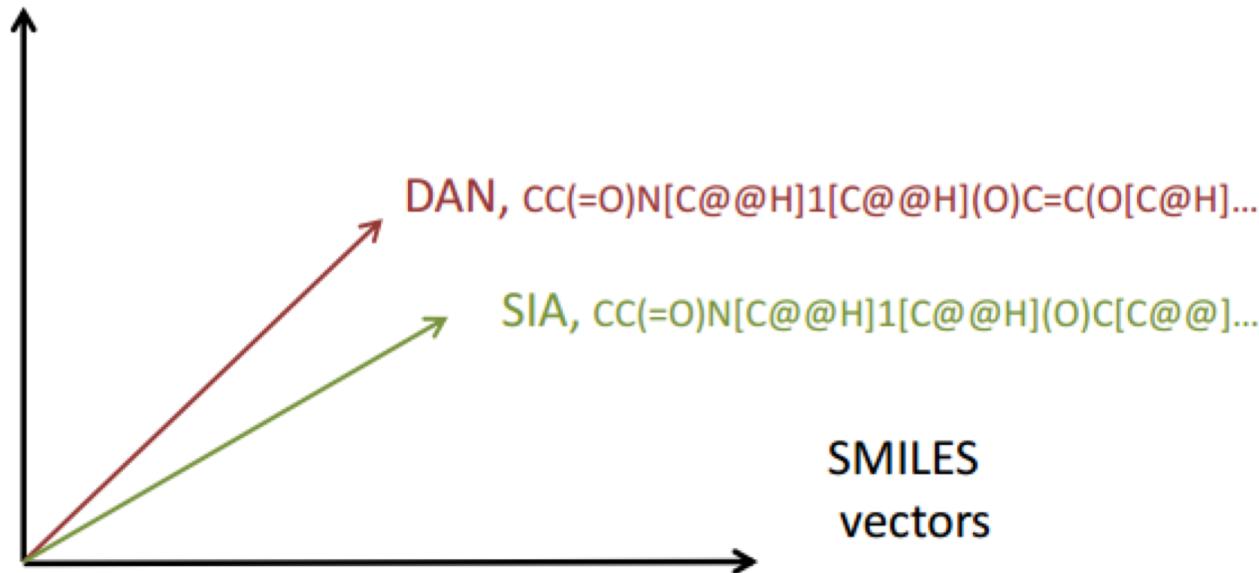
# Öneri: Ligand tabanlı protein temsili



# SMILESVec-tabanlı Protein Temsili

Protein: sialidase

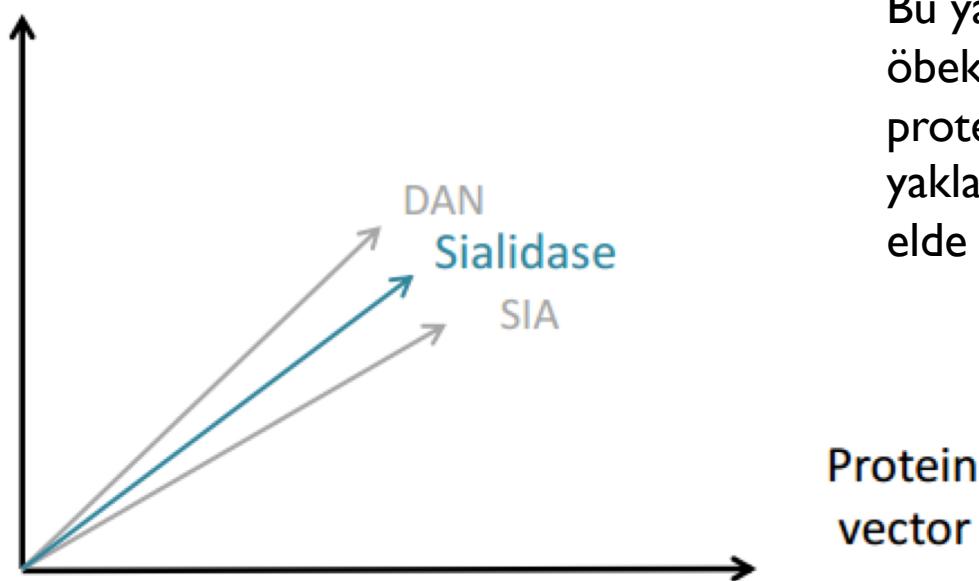
Interacting ligands: DAN, SIA



# SMILESVec-based Protein Representation

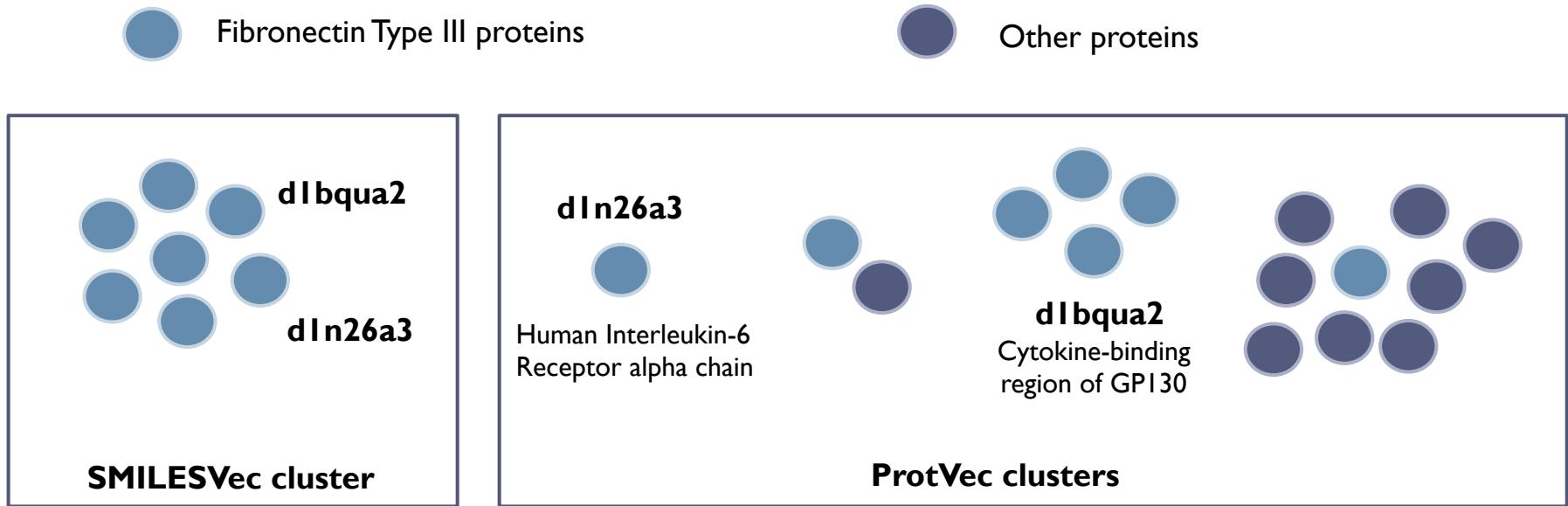
Protein: sialidase

Interacting ligands: DAN, SIA

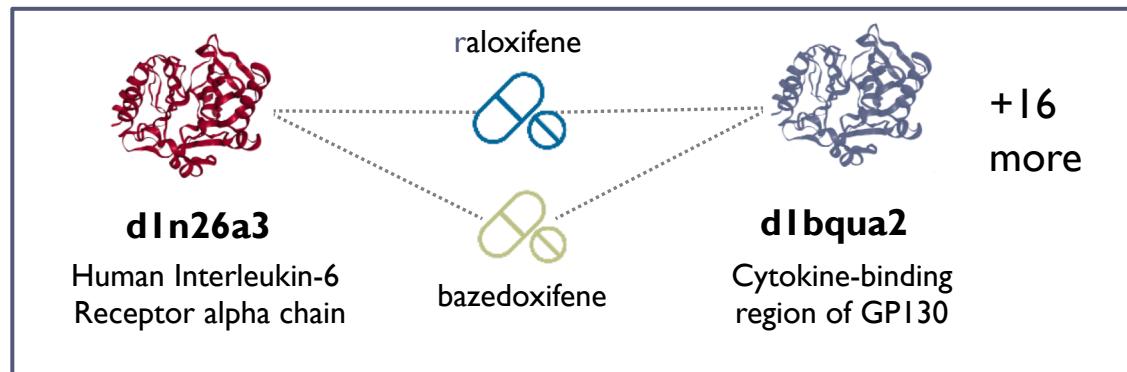
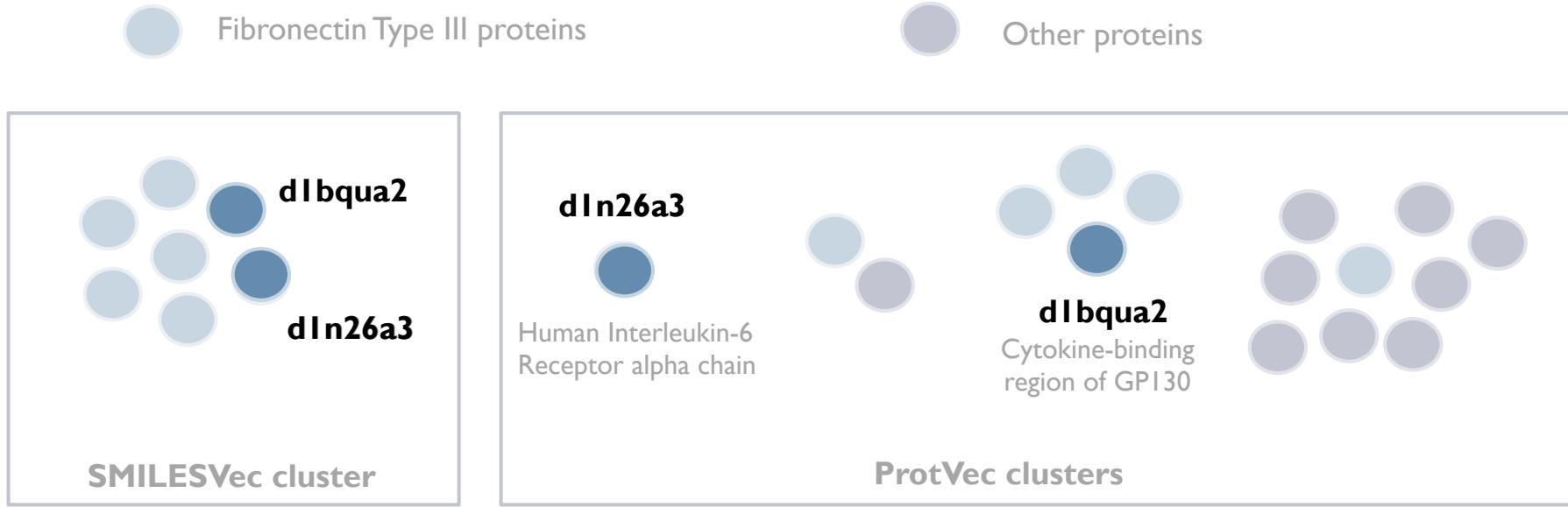


Bu yaklaşımla protein öbekleme probleminde protein dizisini kullanan yaklaşımalarla aynı başarıyı elde ettiğ.

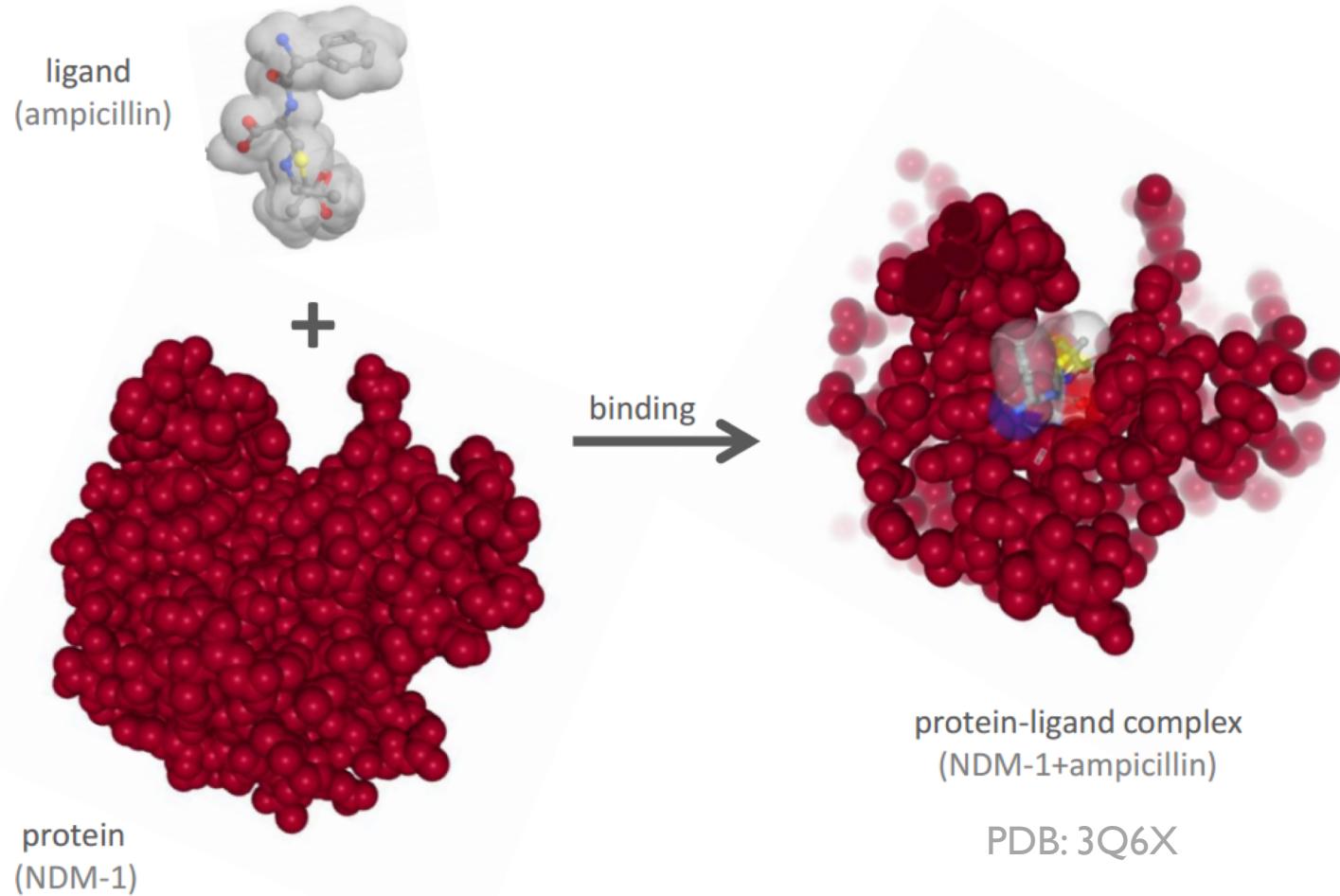
# SMILESVec ve ProtVec ile öbeklemeye örnek



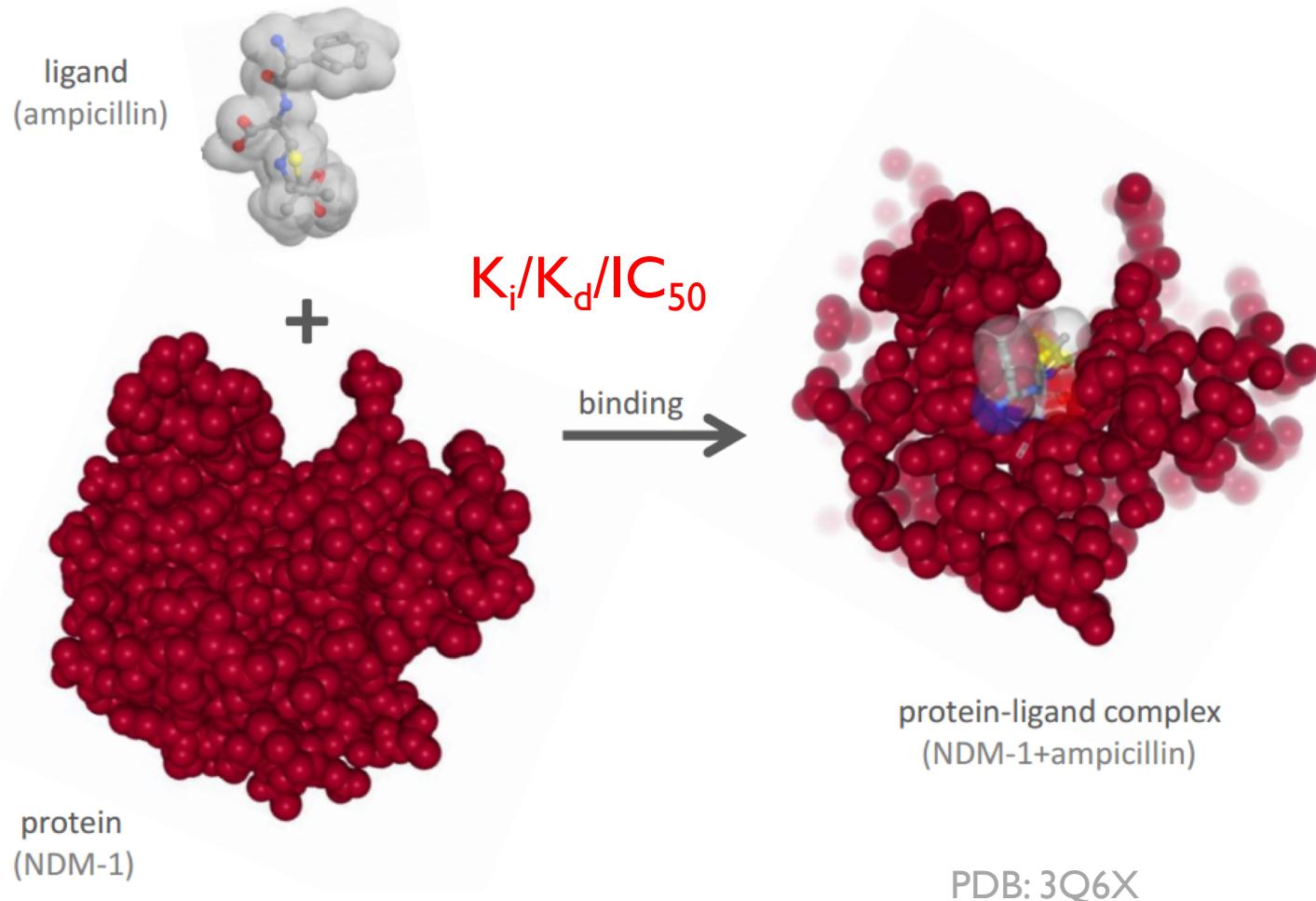
# SMILESVec ve ProtVec ile öbeklemeye örnek



# İlaç-protein etkileşim tahmini



# İlaç-protein bağlanma ilgisi tahmini



# DeepDTA

---

Derin öğrenme algoritması ile kimyasal ve proteinin sadece metin tabanlı dizilerini kullanıyor.



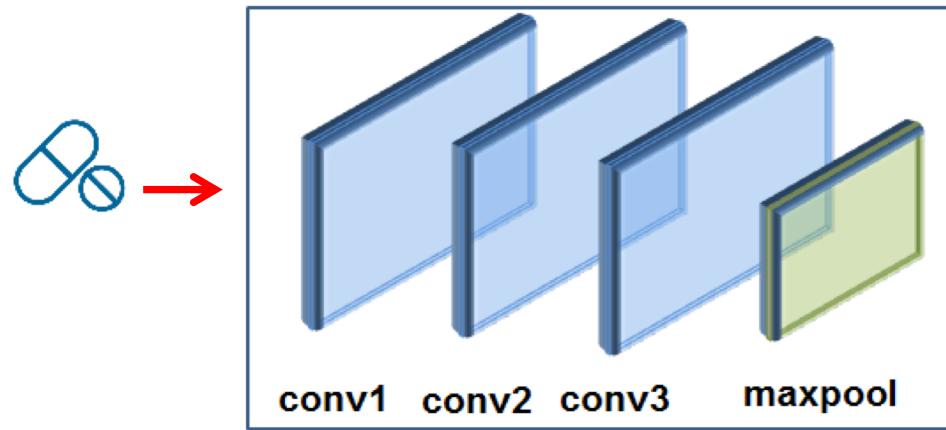
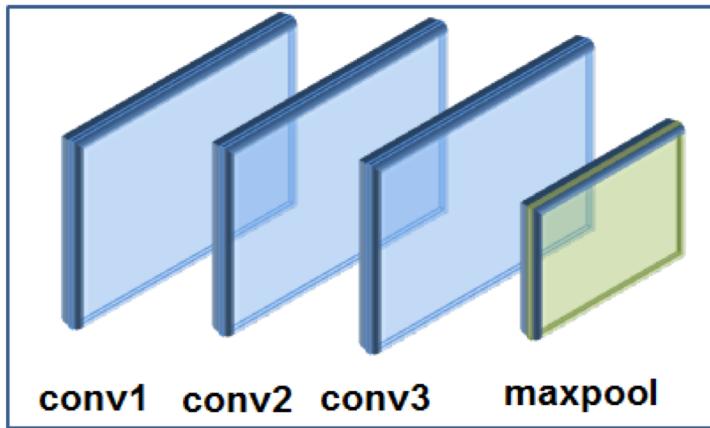
CC1(C(N2C(S1)C(C2=O)NC(=O)...



MELPNIMHPVAKLSTALAAALML..

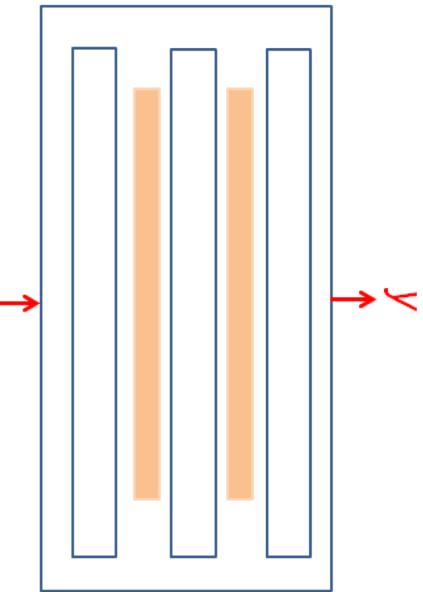


# DeepDTA Modeli



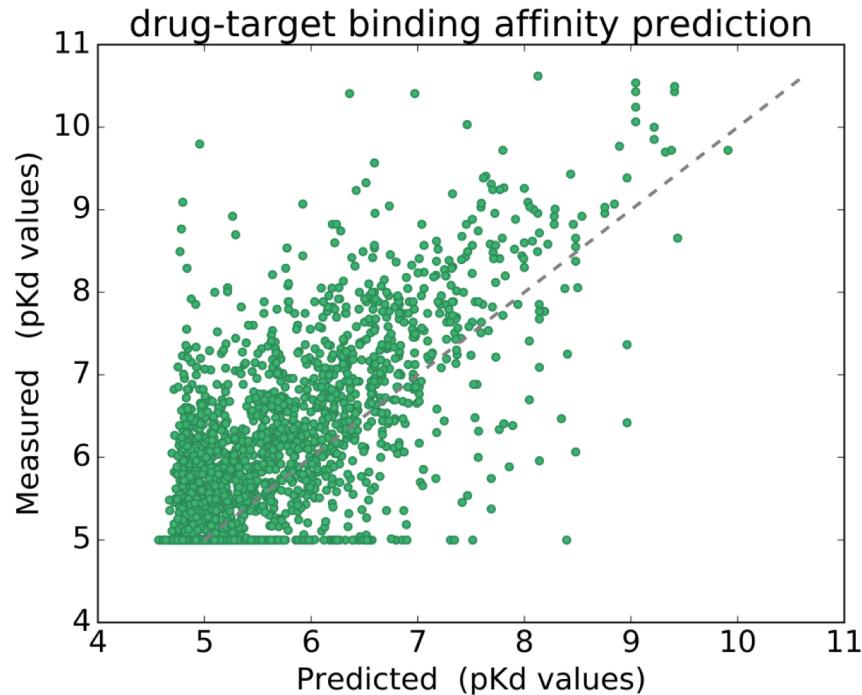
Protein-ilaç temsili için  
Evrişimli Sinir Ağları (CNN)

Protein  
representation  
↑  
Drug  
representation  
↑  
concatenation



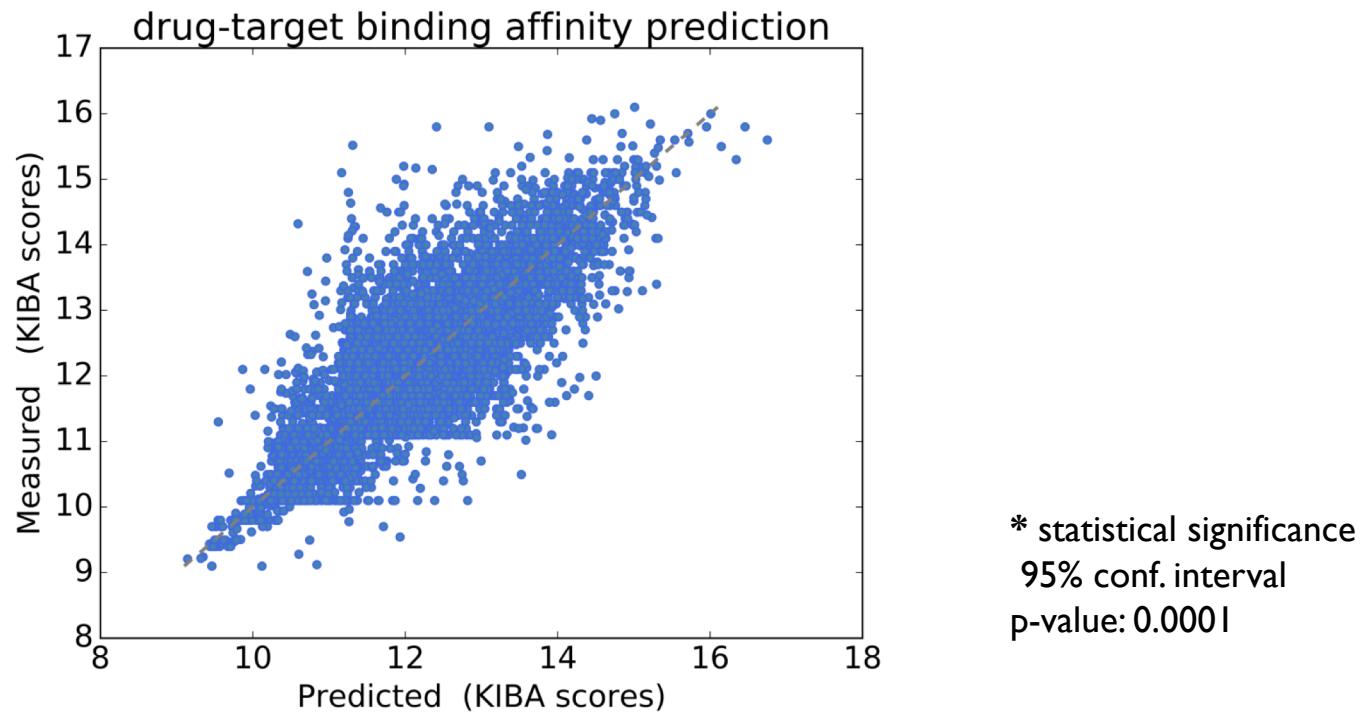
Bağlanması ilgisi tahmini için  
çok katmanlı yapay sinir ağları

# Sonuçlar – Davis veri kümesi



Method	CI	AUPR	MSE
KronRLS (Pahikkala et al., 2014)	0.871	0.661	0.379
SimBoost (He et al., 2017)	0.872	0.709	0.282
<b>DeepDTA</b>	<b>0.878</b>	<b>0.714</b>	<b>0.261</b>

# Sonuçlar – KIBA veri kümesi



Method	CI	AUPR	MSE
KronRLS (Pahikkala et al., 2014)	0.782	0.635	0.411
SimBoost (He et al., 2017)	0.836	0.760	0.222
<b>DeepDTA</b>	<b>0.863*</b>	<b>0.788*</b>	<b>0.194</b>

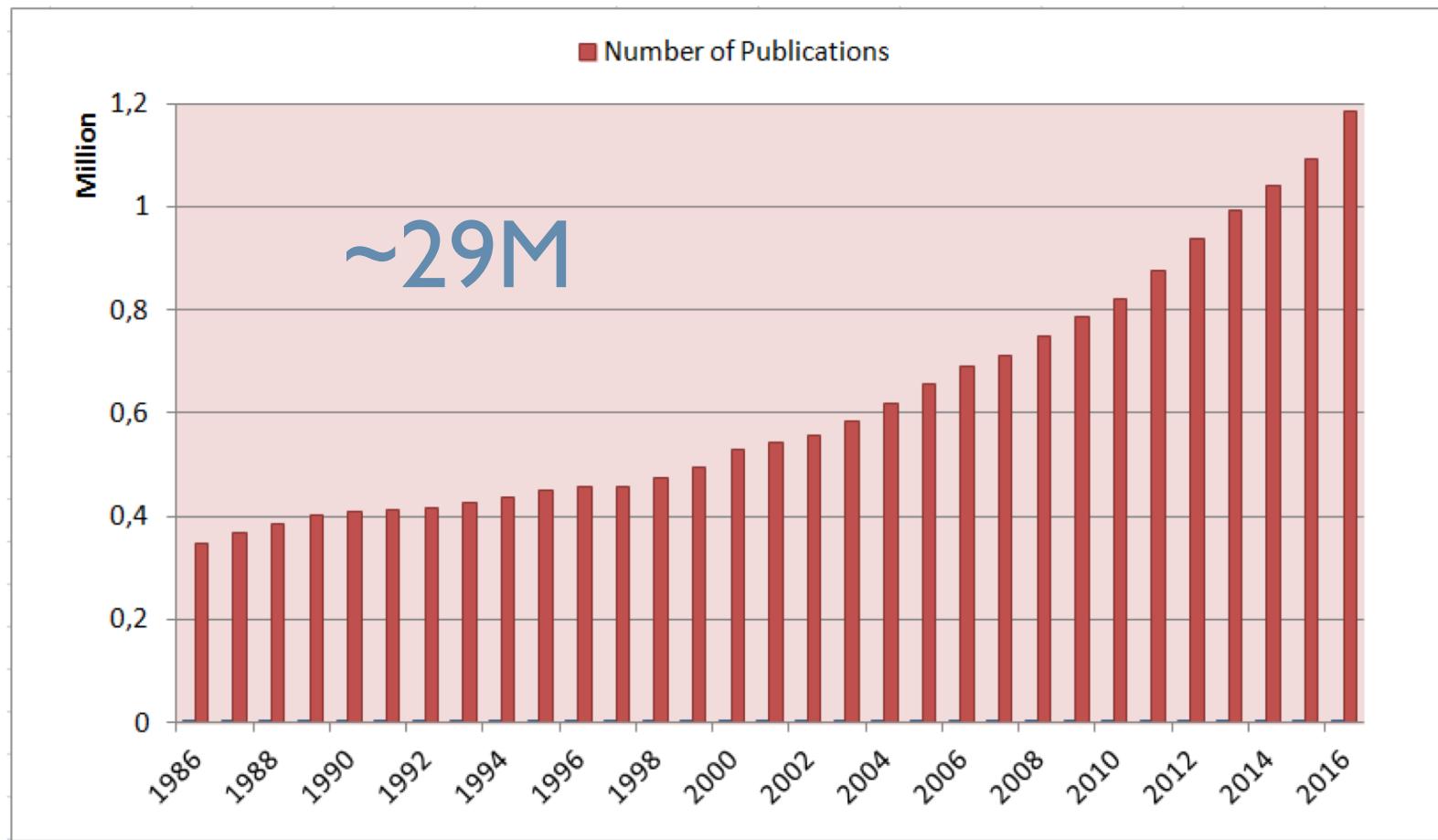
---



# Biyomedikal Metin Madenciliği



# Bilimsel Yayınlarında Artış



# Literatürdeki bilgiye erişim zorluğu

The screenshot shows a PubMed search results page. At the top, the search term 'IFN-gamma OR interferon-gamma' is highlighted with a blue oval. Below the search bar, there are options to 'Create RSS', 'Create alert', and 'Advanced'. The search results are formatted as 'Summary', sorted by 'Most Recent', with 20 items per page. A 'Send to' dropdown menu is visible on the right.

**Best matches for IFN-gamma OR interferon-gamma:**

- [Interferon-gamma: an overview of signals, mechanisms and functions.](#)  
Schroder K et al. J Leukoc Biol. (2004)
- [Cellular responses to interferon-gamma.](#)  
Boehm U et al. Annu Rev Immunol. (1997)
- [\[Interferon \(IFN\) therapy \(recombinant IFN-alpha-2C or recombinant IFN-gamma\) in metastasized hypernephroma\].](#)  
Kuzmits R et al. Acta Med Austriaca. (1985)

[Switch to our new best match sort order](#)

**Search results**  
Items: 1 to 20 of 124125

<< First < Prev Page 1 of 6207 Next > Last >>

[Protective efficacy induced by DNA prime and recombinant protein boost vaccination with Toxoplasma gondii GRA14 in mice.](#)  
1. [Pagheh AS, Sarvi S, Gholami S, Asgarian-Omrani H, Valadan R, Hassannia H, Ahmadpour E, Fasihi-Ramandie M, Dodangeh S, Hosseni-Khah Z, Daryani A. Microb Pathog. 2019 Jun 15:103601. doi: 10.1016/j.micpath.2019.103601. \[Epub ahead of print\]](#)  
PMID: 31212035  
[Similar articles](#)

**Çok zor veya imkansız:** bilim insanların ilgili yayınları takip edebilmesi.

**Manuel oluşturulan veritabanları:** mevcut bilginin çok küçük bir kısmını kapsamaktadır.

# Amaç

---

- ▶ Doğal dil işleme ve yapay öğrenmeye dayanan yöntemlerle önemli bilgilerin otomatik olarak çıkarılması.
- ▶ Gizli bağlantıların tespit edilerek yeni bilimsel hipotezlerin oluşturulması.



# Ne tür bilgiler çıkarabiliriz?

Olumsuzluk  
İlişki  
Türü  
Yön

## Physical and functional interactions between STAT3 and ZIP kinase.

Signal transducer and activator of transcription 3 (STAT3) is a latent cytoplasmic transcription factor that can be activated by cytokines and growth factors. It plays important roles in cell growth, apoptosis and cell transformation, and is constitutively active in a variety of tumor cells. In this study, we provide evidence that zipper-interacting protein kinase (ZIPK) interacts physically with STAT3. ZIPK specifically interacted with STAT3, and did not bind to STAT1, STAT4, STAT5a, STAT5b or STAT6. ZIPK phosphorylated STAT3 on serine 727 (Ser727) and enhanced STAT3 transcriptional activity. Small interfering RNA-mediated reduction of ZIPK expression decreased leukemia inhibitory factor (LIF)- and IL-6-induced STAT3-dependent transcription. Furthermore, LIF- and IL-6-mediated STAT3 activation stimulated ZIPK activity. Taken together, our data suggest that ZIPK interacts with STAT3 within the nucleus to regulate the transcriptional activity of STAT3 via phosphorylation of Ser727.

İlişkiler  
(etkileşimler)

Bağlanma  
yeri

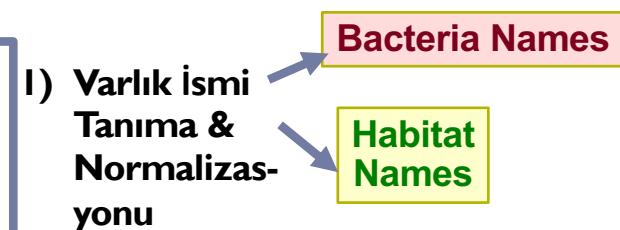
Karmaşık  
olaylar

Spekulatif  
bilgi

Hücresel bölge

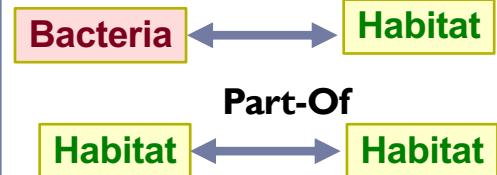
# Ne tür bilgiler çıkarabiliriz?

Bifidobacterium longum NCC2705  
Description  
Bifidobacterium. Representatives of this genus naturally colonize the human gastrointestinal tract (GIT) and are important for establishing and maintaining homeostasis of the intestinal ecosystem to allow for normal digestion. Their presence has been associated with beneficial health effects, such as prevention of diarrhea, amelioration of lactose intolerance, or immunomodulation. The stabilizing effect on GIT microflora is attributed to the capacity of bifidobacteria to produce bacteriocins, which are bacteriostatic agents with a broad spectrum of action, and to their pH-reducing activity. Most of the ~30 known species of bifidobacteria have been isolated from the mammalian GIT, and some from the vaginal and oral cavity. All are obligate anaerobes belonging to the Actinomycetales, branch of Gram-positive bacteria with high GC content that also includes Corynebacteria, Mycobacteria, and Streptomyces.  
Description  
Bifidobacterium longum. This organism is found in adult humans and formula fed infants as a normal component of gut flora.  
Description  
Bifidobacterium longum strain NCC2705. This strain was isolated from infant feces. The genome of this strain is being sequenced for comparative genomics.

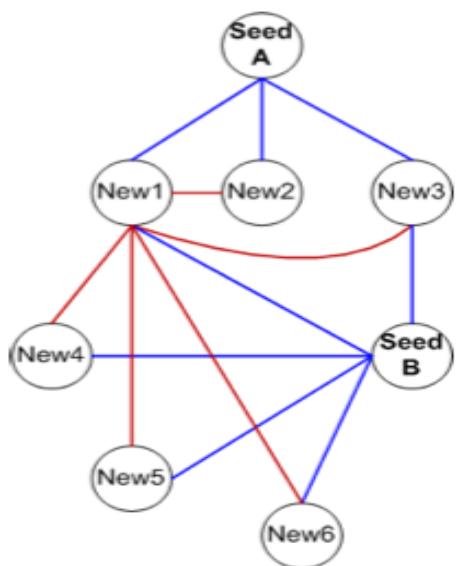


## 2) İlişki Çıkarma

Localization

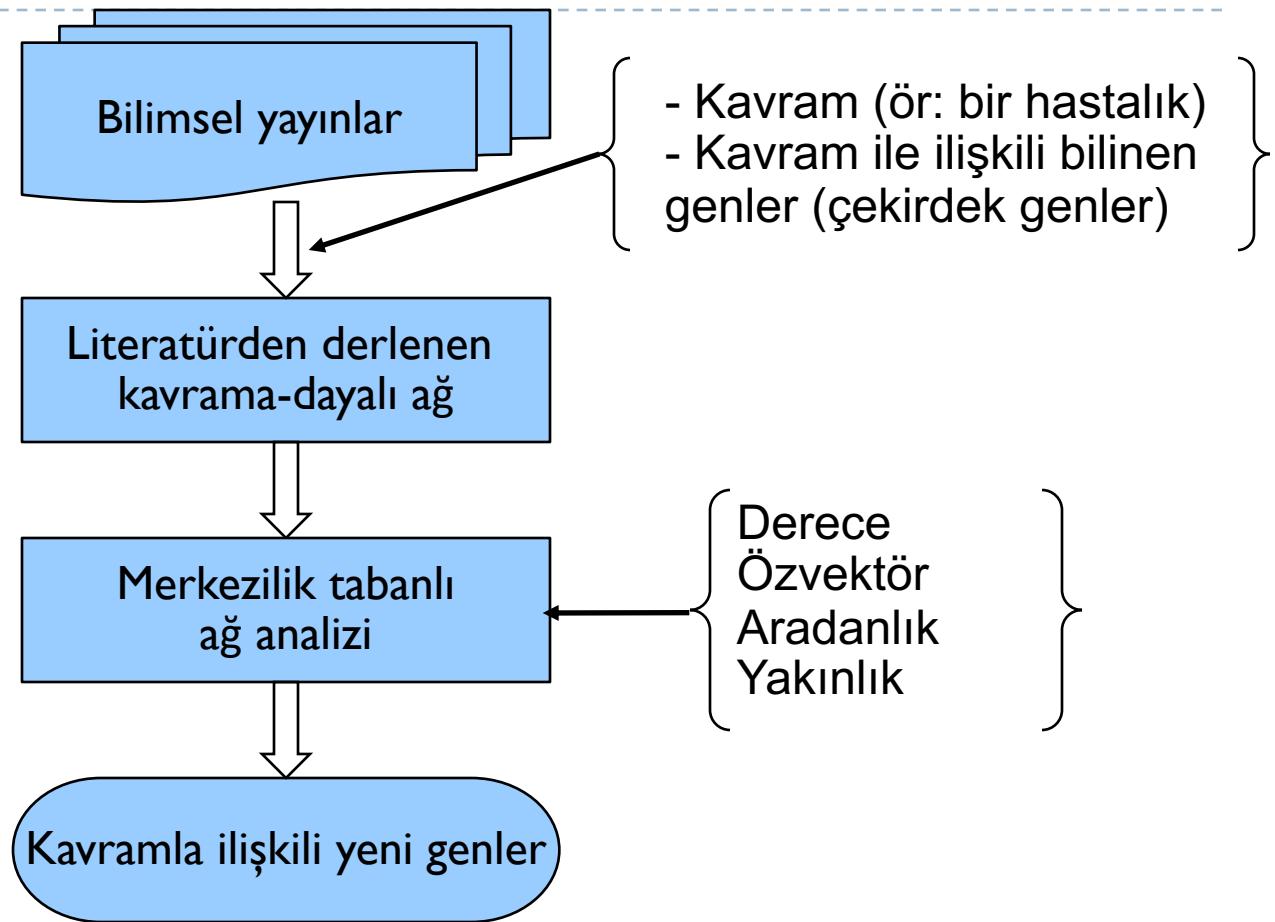


# Metin ve Ağ Madenciliği ile Yeni Bağlantı Tespiti



## Hipotez:

Kavrama bağlı genlerin etkileşim ağındaki önemli genler de bu kavramla ilişkili olabilir.

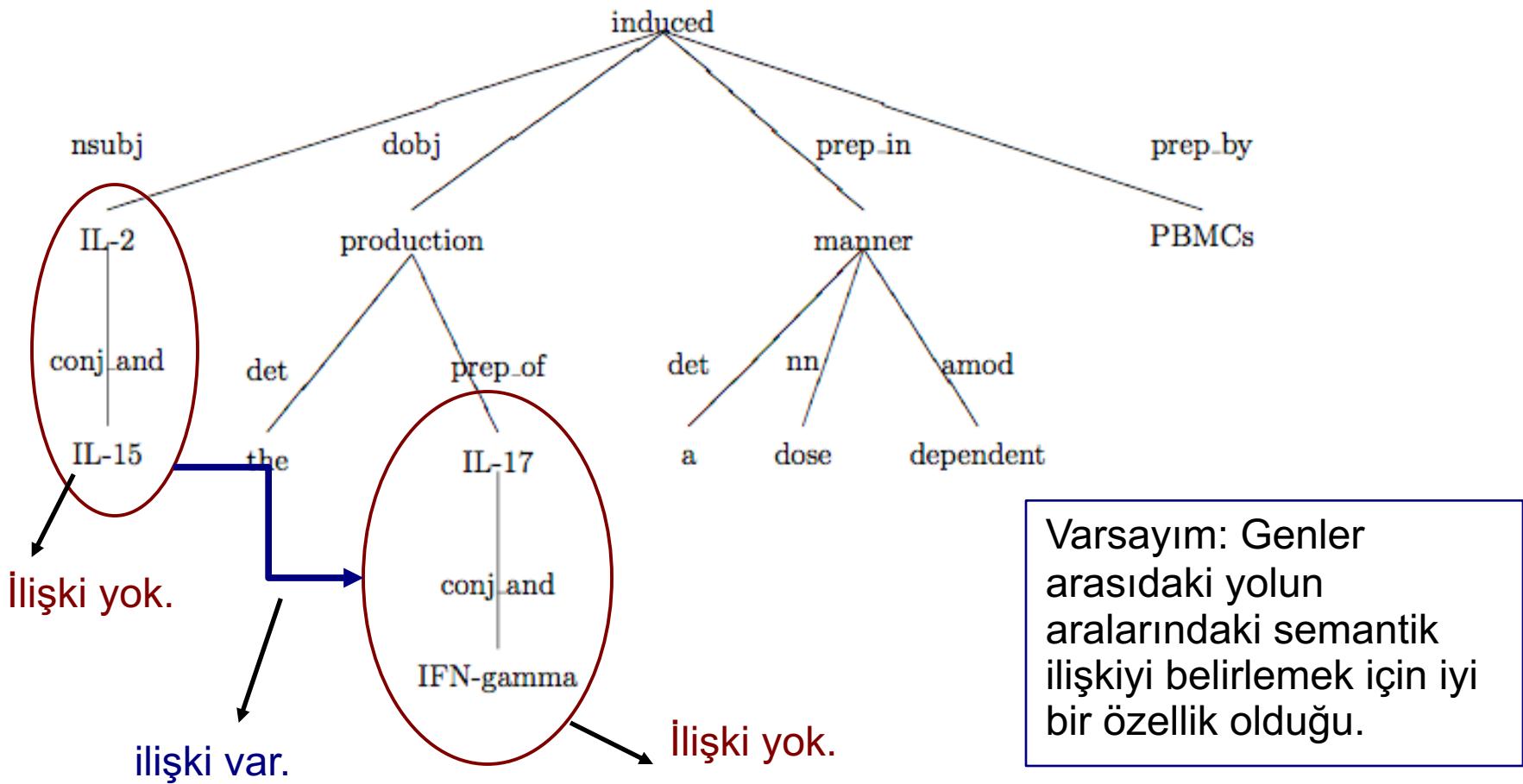


Prostat kanseri, aşıyla bağlı bağımlılık ve ateş ile ilişkili genlerin tespiti için kullanılmıştır.

# Gen-gen etkileşimlerinin tespiti

IL-2 and IL-15 induced the production of IL-17 and IFN-gamma in a dose dependent manner by PBMCs.

(Genia Tagger: 71% F-measure)



► Bağlam ağaçlarının oluşturulması için Stanford Parser kullanıldı (de Marneffe et al., 2006).

# Bağlam Ağacı Yolu Edit Fonksiyonu

- ▶ İlk diziyi ikinciye çevirmek için minimum işlem sayısı (işlemler: bir kelimenin eklenmesi, silinmesi veya değiştirilmesi)
  - ▶ **IL2** – nsubj – induced – dobj – production – prep\_of – **IL-17**
  - ▶ **IL2** – nsubj – induced – dobj – production – prep\_of – IL-17 – conj\_and – **IFN-gamma**
  - ▶ **IL-17** – conj\_and – **IFN-gamma**
    - ▶ Edit mesafesi ( $\text{Yol1} \rightarrow \text{Yol2}$ ) = 2 (2 ekleme)
    - ▶ Edit mesafesi ( $\text{Yol1} \rightarrow \text{Yol3}$ ) = 8 (6 çıkarma + 2 ekleme)

**Benzerlik fonksiyonuna dönüştürme:**

$$EditSim(p_i, p_j) = e^{[-\gamma(EditDist(p_i, p_j))]}$$

- Çekirdek fonksiyon olarak SVM<sup>light</sup> paketine (Joachims, 1999) entegre edildi.

# Değerlendirme

## Veri kümeleri:

Data Set	Sentences	+ Sentences	- Sentences
<b>AIMED<sup>1</sup></b>	4026	951	3075
<b>CB<sup>2</sup></b>	4056	2202	1854

## Sonuçlar:

	Precision	Recall	F-measure
AIMED	77.52	43.51	55.61
CB	85.15	84.79	84.96

**Kesinlik (Precision):** Bulunan ilişkiler içinde doğru olanların oranı.

**Bulma (Recall):** Mevcut ilişkilerden doğru bulunanların oranı.

**F-Ölçütü:** Kesinlik ve bulmanın harmonik ortalaması.

<sup>1</sup> <ftp://ftp.cs.utexas.edu/pub/mooney/bio-data/>

<sup>2</sup> [http://biocreative.sourceforge.net/biocreative\\_2.html](http://biocreative.sourceforge.net/biocreative_2.html)

# Prostat Kanseri Genlerinin Tahmini

OMIM Morbid Map'ten 15 prostat kanseri geni çekirdek gen olarak kullanıldı.

Gene	Description
AR	androgen receptor
BRCA2	breast cancer 2, early onset
MSR1	macrophage scavenger receptor 1
EPHB2	EPH receptor B2
KLF6	Kruppel-like factor 6
MAD1L1	MAD1 mitotic arrest deficient-like 1 (yeast)
HIP1	huntingtin interacting protein 1
CD82	CD82 molecule
ELAC2	elaC homolog 2 (E. coli)
MXI1	MAX interactor 1
PTEN	phosphatase and tensin homolog (mutated in multiple advanced cancers 1)
RNASEL	ribonuclease L (2',5'-oligoisoadenylate synthetase-dependent)
HPC1	hereditary prostate cancer 1
CHEK2	CHK2 checkpoint homolog (S. pombe)
PCAP	predisposing for prostate cancer

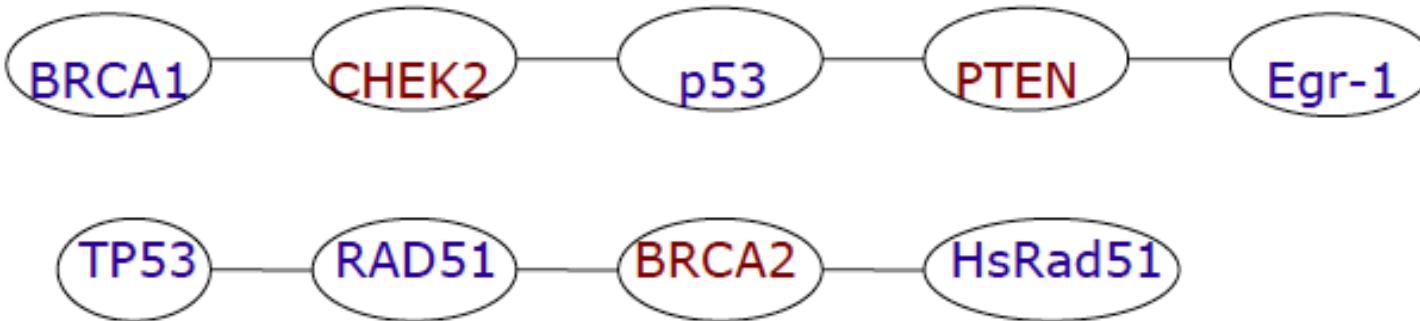
► A. Ozgur, T. Vu, G. Erkan, and D. R. Radev. **Identifying gene-disease associations using centrality on a literature mined gene interaction network.** Bioinformatics, Volume 24, Number 13, pp. i277-i285, 2008.

# Etkileşim Ağının Oluşturulması

- Etkileşim içeren örnek cümleler:

- ▶ **PTEN** is transcriptionally regulated by transcription factors such as p53 and **Egr-1**.
- ▶ In response to DNA damage, the cell-cycle checkpoint kinase **CHEK2** can be activated by ATM kinase to phosphorylate **p53** and **BRCA1**, which are involved in cell-cycle control and apoptosis.
- ▶ The interactions of **RAD51** with **TP53**, RPA and the BRC repeats of **BRCA2** are relatively well understood (see Discussion).
- ▶ The interaction of **BRCA2** with **HsRad51** is significantly more different to both RadA and RecA (Figure 2c).

▶ Oluşturulan ağ:

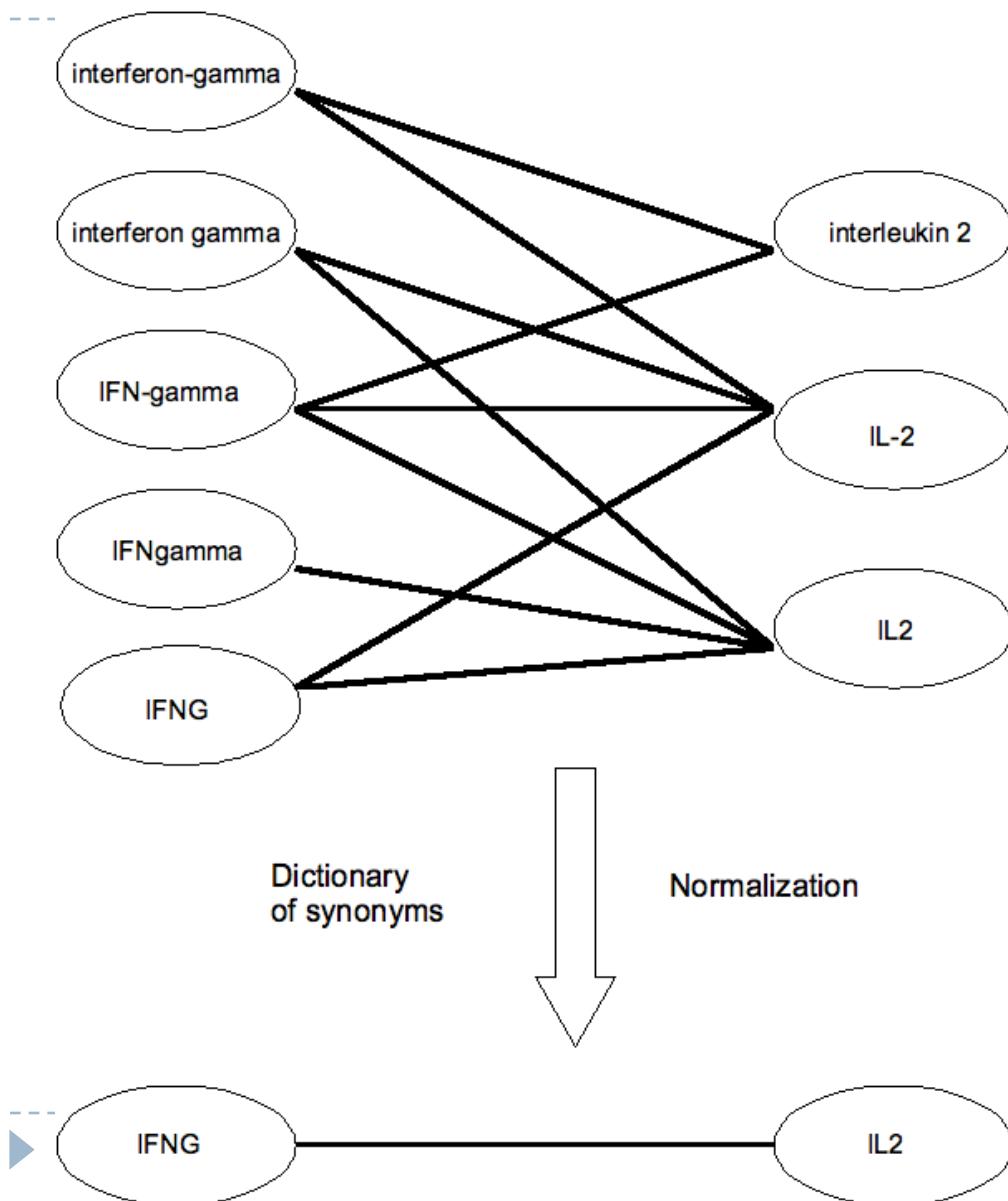


**seed genes:**

CHEK2  
PTEN  
BRCA2



# Gene Adı Normalizasyonu



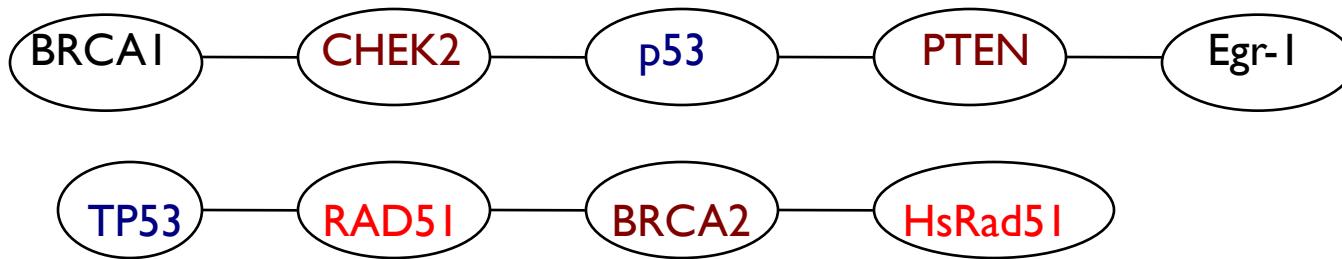
Gen ismi ve eşanlamlıları için  
HUGO Gene Nomenclature  
Committee (HGNC) veritabanı  
sözlük olarak kullanıldı

~28,000 gen girdisi

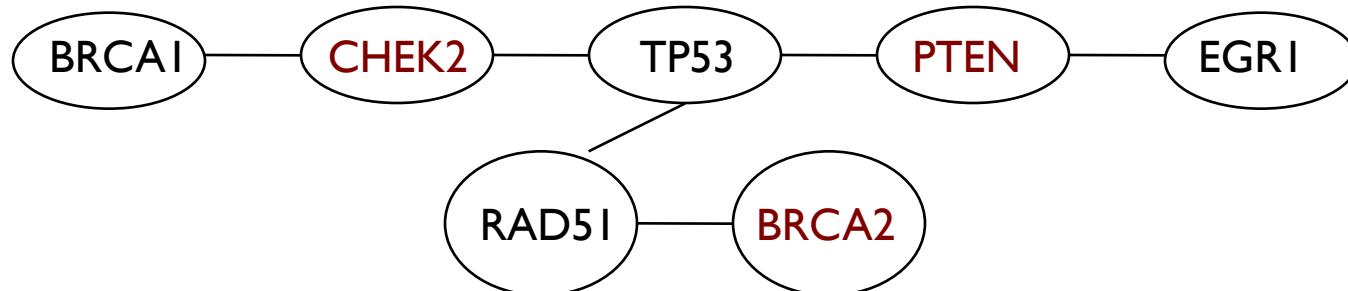
(<http://www.genenames.org/>)

# Etkileşim Ağının Oluşturulması

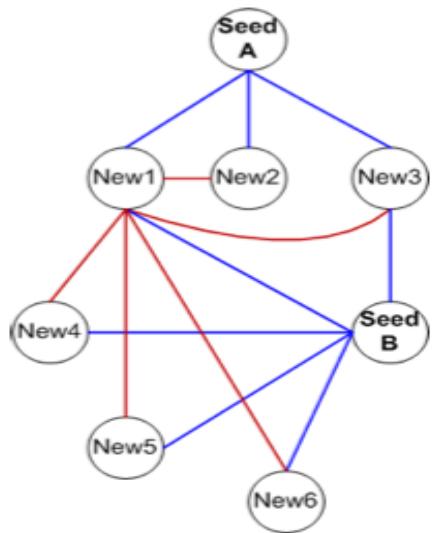
- ▶ Gen ismi normalizasyonundan önce:



- ▶ Gen ismi normalizasyonundan sonra:

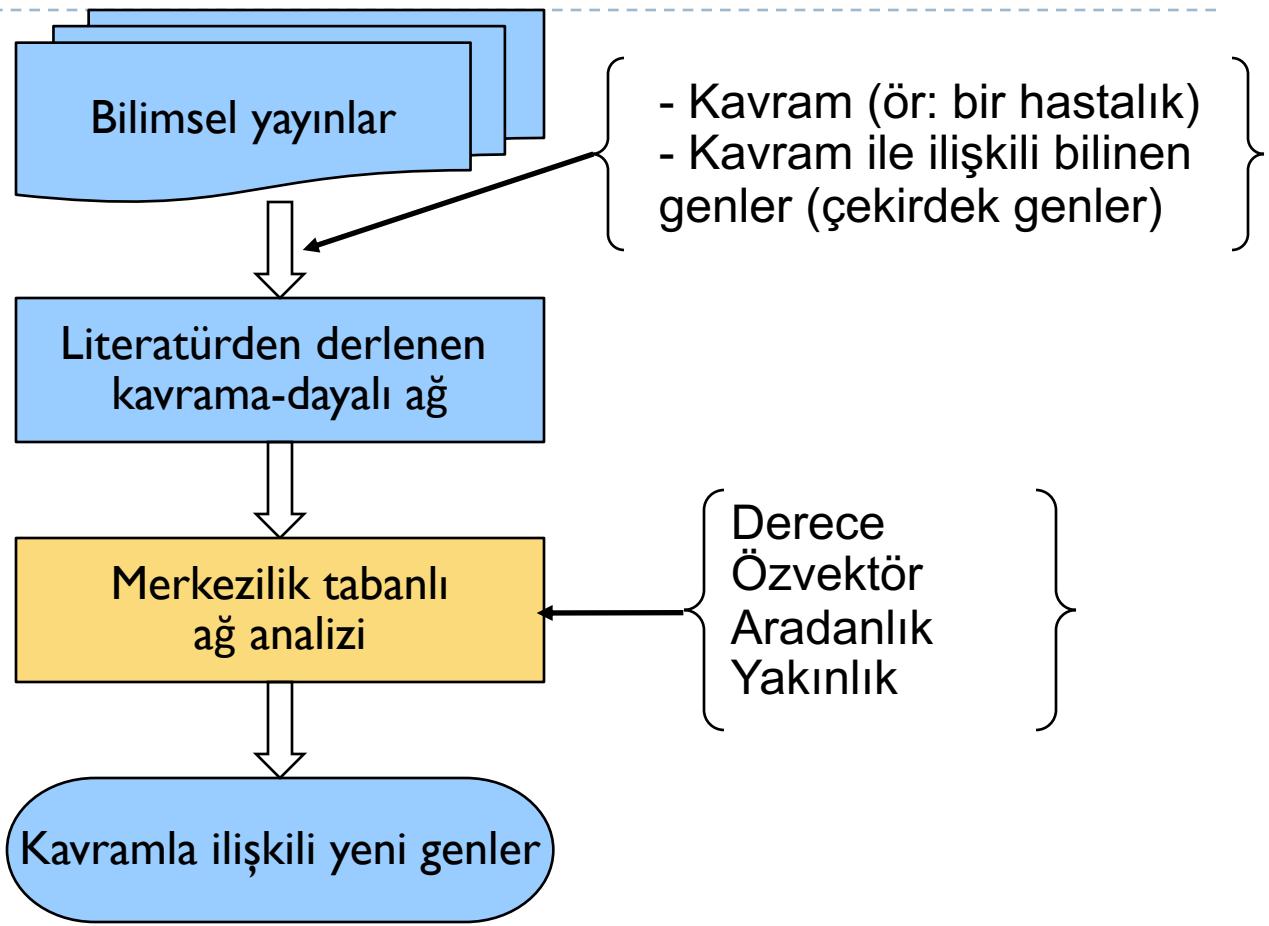


# Genlerin Ağ Merkezilik Ölçütleri ile Sıralanması



## Hipotez:

Kavrama bağlı genlerin etkileşim ağındaki önemli genler de bu kavramla ilişkili olabilir.



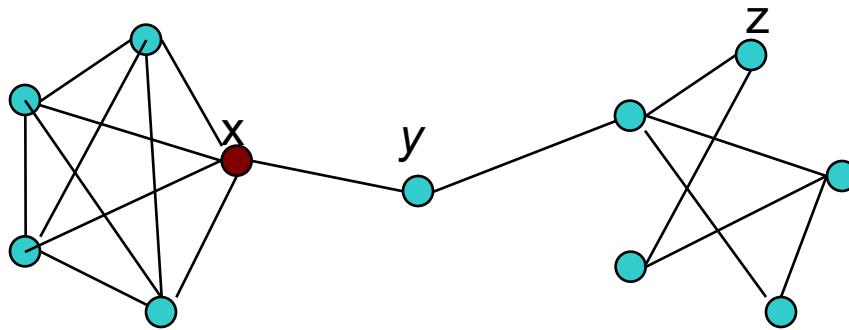
Bir düğümün ağdaki önemi

# Derece Merkezilik Ölçütü

- ▶ Bir düğümün bağlı olduğu düğüm sayısı.

A: Komşuluk matrisi

$$k_i = \sum_{j=1}^n A_{ij}$$



- ▶ Bir düğümün ne kadar çok komşusu varsa, o kadar önemlidir.
- ▶ x düğümünün derece merkezilik ölçütü 5, y'ninki ise 2.

# Özvektör Merkezilik Ölçütü

---

- ▶ Bir düğümün komşularının merkezilik ölçütlerinin toplamıyla orantılıdır.

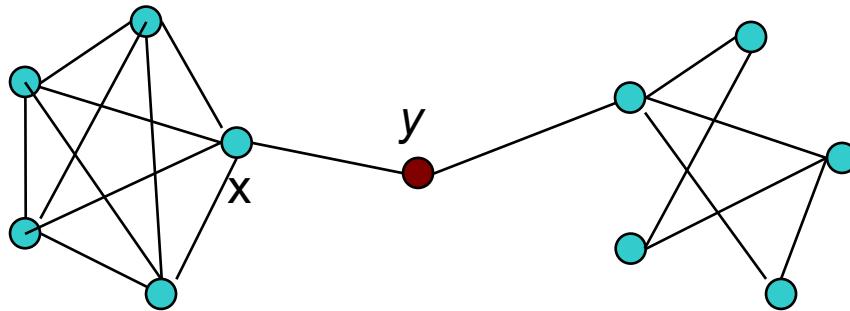
$$x_i = \lambda^{-1} \sum_{j=1}^n A_{ij} x_j$$

- ▶ Matris gösteriminde:  $\lambda \mathbf{x} = \mathbf{Ax}$ 
  - ◆  $\lambda$   $\mathbf{A}$  matrisinin en büyük özdeğeri,  $\mathbf{x}$  de ilgili özvektördür.
- ▶ Her komşu bir düğümün merkeziliğine eşit oranda katkı sağlamamaktadır.
- ▶ Sosyal ağlarda “prestij” olarak ifade edilmektedir.
  - ▶ Bir kişinin prestiji sadece kaç arkadaşı olduğuna değil, bu arkadaşların kim (ne kadar prestijli) olduğuna da bağlıdır.



# Yakınlık Merkezilik Ölçütü

- I/(düğümün diğer düğümlere olan uzaklıklarının toplamı)

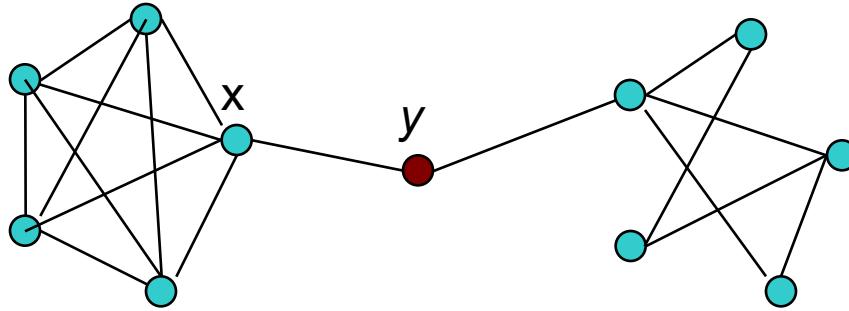


- Bir düğüm diğer düğümlere ne kadar yakınsa, o kadar önemlidir.



# Aradanlık Merkezilik Ölçütü

- Bir  $i$  düğümü için,  $i$  düğümünün üzerinden geçen en kısa yolların sayısının tüm en kısa yollara oranı.



- Bir düğüm ne kadar çok en kısa yol üzerinde yer alırsa, o kadar önemlidir.



# En üst sıradaki 20 Gen

Gene	Degree	Eigenvector	Closeness	Betweenness	Evidence
TP53	+	+	+	+	PGDB
BRCA1	+	+	+	+	PGDB
EREG	+	+	+	+	None
AKT1	+	+	+	+	PGDB
MAPK1	+	+	+	+	Literature ( <a href="#">Hao et al., 2007</a> ; <a href="#">Sarfaraz et al., 2006</a> )
TNF	+	+	+	+	PGDB
CCND1	+	+	+	+	PGDB
MYC	+	+	+	+	PGDB
APC	+	+	-	-	PGDB
CDKN1B	+	+	+	-	PGDB
MAPK8	+	+	+	+	PGDB
NR3C1	-	+	+	-	Literature ( <a href="#">Wei et al., 2007</a> )
VEGFA	+	+	+	-	PGDB
MDM2	+	+	+	-	KEGG and Literature ( <a href="#">Wang et al., 2003</a> ; <a href="#">Zhang et al., 2003</a> )
POLD1	-	-	+	+	None
SNORA62	-	-	+	+	None
CNTN2	-	-	-	+	None
PPA1	-	-	-	+	None
TMEM37	-	-	+	-	None
FZR1	-	-	+	-	PGDB
SSSCA1	-	-	+	-	None
BCL2	+	-	-	-	PGDB
INS	+	-	-	-	KEGG and Literature ( <a href="#">Ho et al., 2003</a> )

12 gen: Prostate Gene DataBase (PGDB)

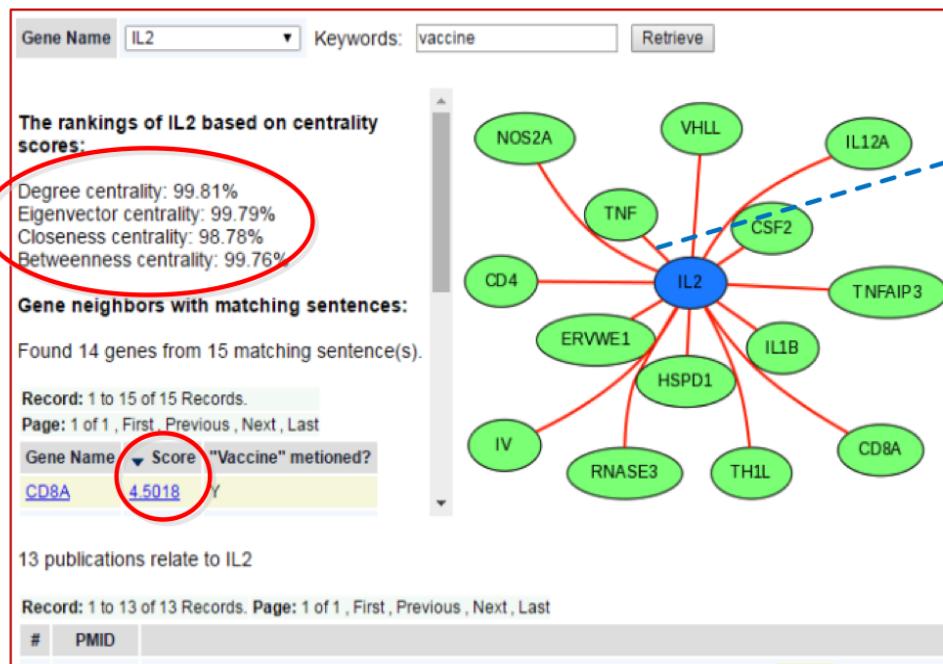
2 gen: KEGG pathway for prostate cancer ve literatür (MDM2 and INS)

2 gen: literatür (NR3C1 and MAPK1)

7 gen: Olumlu veya olumsuz kanıt bulunamadı.

# IGNET: Integrated Gene Network

<http://ignet.hegroup.org>



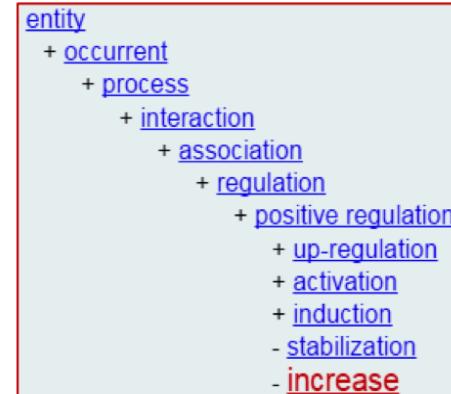
(A)

Gene1: IL2    Gene2: TNF    Keywords:    Search

Found 329 record(s).

Record: 1 to 50 of 329 Records. Page: 1 of 7, First, Previous, Next, Last

PubMed	Gene1	Gene2	Match1	Match2	Score	"Vaccine" mentioned	Sentence	INO Interaction
<a href="#">2422147</a>	IL2	TNF	il 2	TNF	0.86995216	0	Injection of recombinant IFN-gamma with OK-432 or of IFN-alpha/beta, recombinant IFN-beta, recombinant IFN-alpha A/D or recombinant <b>IL2</b> six hours before OK-432 <b>enhanced TNF</b> production about 10-fold, which indicated priming actions of these compounds in <b>TNF</b> production.	increase
<a href="#">3029221</a>	IL2	TNF	interleukin 2	TNF	0.35022627	0	Comparison of <b>TNF</b> receptor <b>expression</b> with that of high affinity <b>interleukin 2</b> (IL-2) and interferon-gamma (IFN-gamma) receptors, respectively, revealed similarities to IL 2-receptor <b>expression</b> with respect to kinetics of induction.	gene expression



(B)

(C)

## Dynamic Ignet Search Program

Keywords Search (through PubMed search engine)

Keywords: bipolar disorder

Degree centrality			Eigenvector centrality		
#	Gene	Score	#	Gene	Score
1	<a href="#">GSK3B</a>	0.0645	1	<a href="#">GSK3B</a>	0.0241
2	<a href="#">AKT1</a>	0.0553	2	<a href="#">AKT1</a>	0.0214
3	<a href="#">PER1</a>	0.0415	3	<a href="#">BDNF</a>	0.0188
4	<a href="#">NFI3L</a>	0.0415	4	<a href="#">DISC1</a>	0.0161
5	<a href="#">BDNF</a>	0.0415	5	<a href="#">HSPA5</a>	0.0123

Keywords: bipolar disorder.

Found 199 gene pairs.

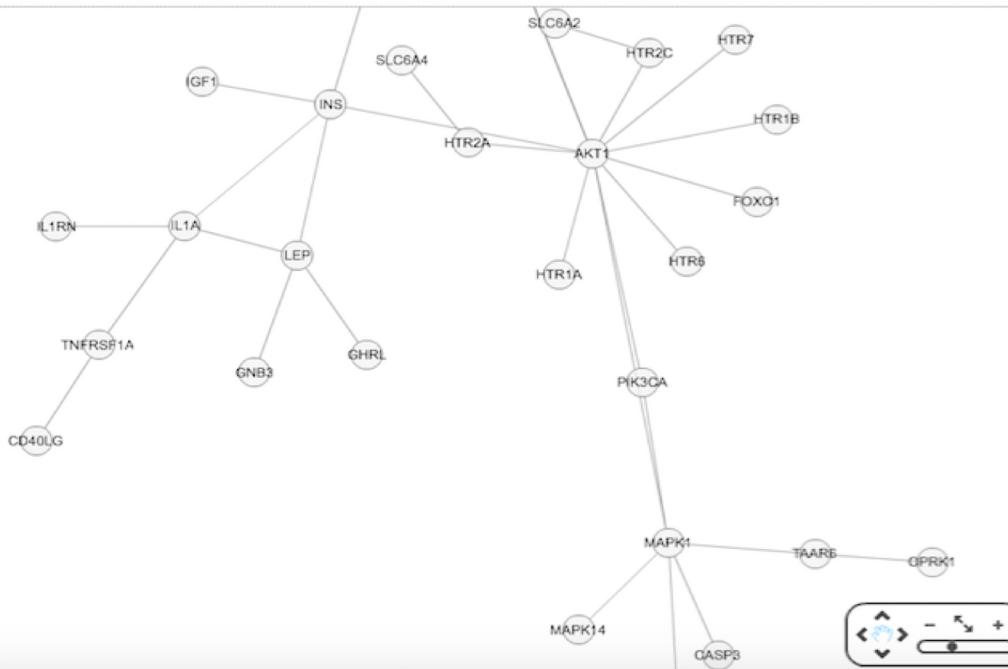
[Calculate centrality scores.](#)

[Show network in Cytoscape Web.](#)

#	Gene 1	Gene 2	Number of hits
1	<a href="#">ABCA1</a>	<a href="#">NR1H3</a>	<a href="#">1 hits</a>
2	<a href="#">ACE</a>	<a href="#">AGT</a>	<a href="#">2 hits</a>
3	<a href="#">ACTB</a>	<a href="#">RELN</a>	<a href="#">1 hits</a>
4	<a href="#">ADA</a>	<a href="#">ADARB1</a>	<a href="#">1 hits</a>
5	<a href="#">ADCY7</a>	<a href="#">SUCLG1</a>	<a href="#">1 hits</a>

Found 199 gene pairs. Below are the network shown in Cytoscape Web.

[Download network in graphml format.](#)



# Yerel Olmayan Bağlam Bilgisi Çıkarma

Olumsuzluk

İlişki  
Türü

Yön

Deney Türü  
Organizma

## Physical and functional interactions between STAT3 and ZIP kinase.

Signal transducer and activator of transcription 3 (STAT3) is a latent cytoplasmic transcription factor that can be activated by cytokines and growth factors. It plays important roles in cell growth, apoptosis and cell transformation, and is constitutively active in a variety of tumor cells. In this study, we provide evidence that zipper-interacting protein kinase (ZIPK) interacts physically with STAT3. ZIPK specifically interacted with STAT3, and did not bind to STAT1, STAT4, STAT5a, STAT5b or STAT6. ZIPK phosphorylated STAT3 on serine 727 (Ser727) and enhanced STAT3 transcriptional activity. Small interfering RNA-mediated reduction of ZIPK expression decreased leukemia inhibitory factor (LIF)- and IL-6-induced STAT3-dependent transcription. Furthermore, LIF- and IL-6-mediated STAT3 activation stimulated ZIPK activity. Taken together, our data suggest that ZIPK interacts with STAT3 within the nucleus to regulate the transcriptional activity of STAT3 via phosphorylation of Ser727.

İlişkiler  
(etkileşimler)

Bağlanma  
yeri

Karmaşık  
olaylar

Spekulatif  
bilgi

Hücresel bölge

Makalenin  
tüm metni

---

# Deneysel Yöntemlerin Açıklandığı Pasajları Bulma

İşbirliği: Ferhat Aydın ve Zehra Melce Hüsünbeyi

Aydın F., Hüsünbeyi Z.M., Özgür A. (2017). **Automatic query generation using word embeddings for retrieving passages describing experimental methods.** *Database*, 2017.

---



# Birden çok deneysel yöntemin anlatıldığı örnek bir paragraf

---

To identify novel interactors of TBK1, we generated stable RAW264.7 cell lines that express a GS-TAP-tagged version of TBK1 and purified the TBK1 protein complex by tandem affinity purification (Figure 1A and B) (Burckstummer et al. 2006). Mass spectrometry analysis identified TBK1 along with the core complex components TRF family member-associated NF-kappa-B activator (TANK; 17 peptides), TBK-binding protein 1 (TBKBP1, also referred to as SINTBAD; 10 peptides) and TBKBP2 (also referred to as NAPI or AZI2; 15 peptides), indicating that this native purification was efficient and in agreement with previously published data on the TBK1 core complex (Bouwmeester et al. 2004). In addition, we identified the DEAD-box helicase DDX3X (RefSeq NP\_034158) with five peptides. Immunoprecipitation experiments using tagged TBK1 suggested that its interaction with DDX3X and the transcription factor IRF3 are significantly weaker than the interaction between TBK1 and TANK and therefore not detected by coimmunoprecipitation under stringent conditions (Supplementary Figure 1A). Likewise, we could not detect any association of DDX3X with IRF3 (Supplementary Figure 1B).



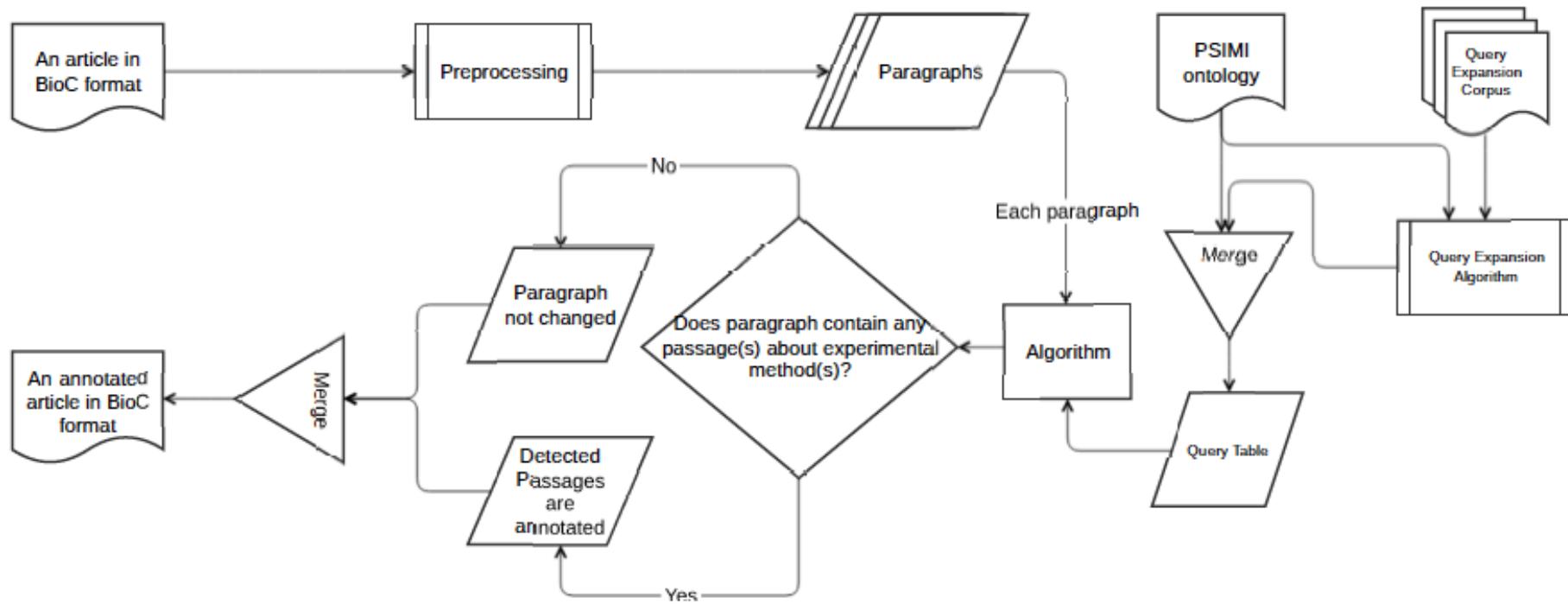
# Manuel Etiketleme

---

```
<passage>
  <infon key="type">paragraph</infon>
  <offset>17753</offset>
  <text>Infection of eukaryotic cells with large DNA viruses often results in extensive interactions of viral gene products with macromolecular pathways of the host cell. By using the yeast two-hybrid system, we identified cellular hnRNP-K as an interacting protein with ASFV early protein p30. This interaction was further confirmed by an in vitro GST-fusion pull-down assay, using either p30 obtained from baculovirus system or ASFV infected cell extracts.</text>
  <annotation id="10">
    <infon key="type">ExperimentalMethod</infon>
    <infon key="PSIMI">0018</infon>
    <location offset="17916" length="123"/>
    <text>By using the yeast two-hybrid system, we identified cellular hnRNP-K as an interacting protein with ASFV early protein p30.</text>
  </annotation>
  <annotation id="11">
    <infon key="type">ExperimentalMethod</infon>
    <infon key="PSIMI">0096</infon>
    <location offset="18040" length="163"/>
    <text>This interaction was further confirmed by an in vitro GST-fusion pull-down assay, using either p30 obtained from baculovirus system or ASFV infected cell extracts.</text>
  </annotation>
</passage>
```



# Sistemin genel iş akışı



# Sorgu oluşturma

```
{  
    name = "0018"  
    synonym = [ two hybrid, two-hybrid, yeast two hybrid, 2 hybrid, 2-hybrid, 2h, y2h,  
                classical two hybrid, gal4 transcription regeneration, yeast two-hybrid ]  
    tier 1 = [ yeast, hybrid, y-2h ]  
    tier 2 = [ bait, cdna, gal4, gal, galactosidase ]  
}  
  
{  
    name = "0019"  
    synonym = [ co-immunoprecipitation, coimmunoprecipitation, co-ip, coip, immunoprecipitation,  
                anti bait coip, anti bait coimmunoprecipitation, anti tag coip,  
                anti tag coimmunoprecipitation ]  
    tier 1 = [ immunoprecipitated, immunoprecipitates, precipitated, precipitate, precipitates,  
                co-precipitated, ip ]  
    tier 2 = [ antibody, antibodies, tag, tagged, bait, immunoblotted]  
}  
  
{  
    name = "0096"  
    synonym = [ pull down, affinity capture, pulldown, pull-down, pulled-down, pulled down ]  
    tier 1 = [ pull, pulled, down, affinity, capture ]  
    tier 2 = [ gst, appl1, rab5, gel, glutathione ]  
}
```



İsim ve eşanlamlılar PSI-MI ontolojisinden alındı. Tier 1 and Tier 2 terimleri tf.rf ile bulundu.

# Sorgu cevaplanması

---

- To assess whether acetylation of any of these lysines affected Mediator's ability to interact with the Histone H4 or H3 tails we performed the Mediator ***pull down*** assay with Ac-(K5,8,12,16)-H4 and Ac-(K9,14)-H3 peptides. The Ac-(K9,14)-H3 peptide showed no difference in ***affinity*** from the unmodified H3-N peptide (Figure 5A), even at reduced concentrations of peptide and Mediator.
  
- Cell lysates were subjected to an anti-Flag ***immunoprecipitation (IP)***, and ***coprecipitating*** STRAP was detected by immunoblotting (Blot) with anti-HA ***antibodies*** (top section). In the middle section, total lysates were ***immunoprecipitated*** using anti-HA ***antibodies*** and then immunoblotted with anti-Flag ***antibodies***. To confirm expression of Smads, aliquots of total cell lysates were ***immunoblotted*** with anti-Flag ***antibodies*** (bottom section).



# Sonuçlar

---

	Precision	Recall	F-measure
<b>baseline</b>	0.424	0.418	0.421
<b>baseline.genia.ino</b>	0.484	0.413	0.446
<b>tf.rf.f7s7</b>	0.120	0.508	0.194
<b>tf.rf.f7s7.genia.ino</b>	0.133	0.502	0.211
<b>tf.rf.f10s10</b>	0.068	0.512	0.119
<b>tf.rf.f10s10.genia.ino</b>	0.074	0.507	0.129
<b>tf.rf.manual</b>	0.315	0.508	0.389
<b>tf.rf.manual.genia.ino</b>	0.357	0.503	0.418
<b>word2vec</b>	0.321	0.618	0.422
<b>word2vec.genia.ino</b>	0.362	0.606	0.453
<b>IAA</b>	0.787	0.937	0.856





# Neurobilim Alanında Metin Madenciliği

İşbirliği: Erinç Gökdeniz ve Reşit Canbeyli

E. Gokdeniz, A. Ozgur, R. Canbeyli. **Automated Neuroanatomical Relation Extraction: A Linguistically Motivated Approach with a PVT Connectivity Graph Case Study.** *Frontiers in Neuroinformatics*, 10:39, 2016.

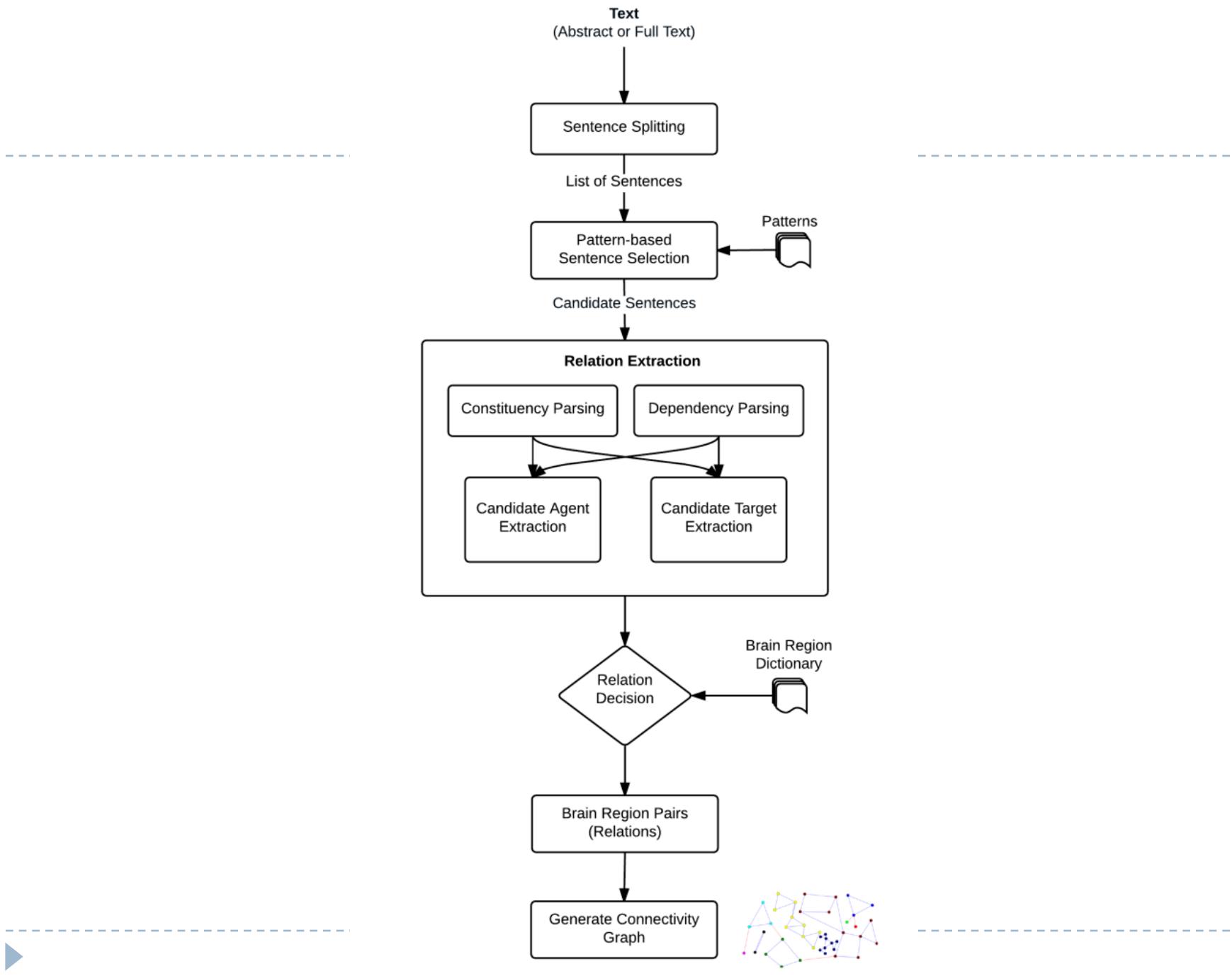
# Problem Tanimi

smaller injection of FB in the medial caudate in OM39 produced a similar pattern of label in Pa, with little or no label in Pt, while an injection of DY in the lateral caudate in the same case labeled virtually no neurons in the midline nuclei (not illustrated, Table 1).

Pa projects strongly to the accumbens nucleus  
strongly to the nucleus. Pt also projects to the shell) and more than the shell.

Amyg projects to the rostromedial caudate nucleus. The magnocellular basal nucleus of the amygdala, with involvement of the accessory basal and lateral nuclei. This experiment had substantial numbers of retrogradely labeled neurons in Pa, Pt, and parts of the Cdc (Fig. 6A-E, Table 1). Lower numbers of labeled neurons were seen in the Clc, Re, Pcn, and Pf. A similar, but smaller injection of CtB into the amygdala in

Hsu et al. (2009)



# İlişki ifade eden kelimeler

---

*“An anterograde tracer injection into the dorsal midline thalamus revealed strong projections to the accumbens nucleus.” [1]*

*“For example, the pPVT was found to be distinctively innervated by the anterior most aspect of the prelimbic cortex and the agranular portions of the posterior insular cortex .” [2]*

# Beyin bölgeleri sözlüğü

Brain Region	Acronym	Synonym
parietal lobe	PL	parietal cortex, parietal region, lobus parietalis
suprachiasmatic nucleus	SCN	suprachiasmatic nuclei
cingulate gyrus	CgG	cingular gyrus, cingulate area, cingulate region, gyri cinguli, gyrus cinguli
subthalamus	SbTh	subthalamic region, ventral thalamus, thalamus ventralis
superior frontal gyrus	SFG	marginal gyrus, superior frontal convolution, gyrus frontalis superior
parabrachial nucleus	-	parabrachial nuclei,parabrachial
paracentral nucleus	PC	paracentral thalamic nucleus, nucleus paracentralis, paracentral nucleus of the thalamus,paracentral
central medial nucleus	CM	central medial thalamic nucleus, nucleus centralis medialis, centralis medialis,central medial nucleus of the thalamus,central medial

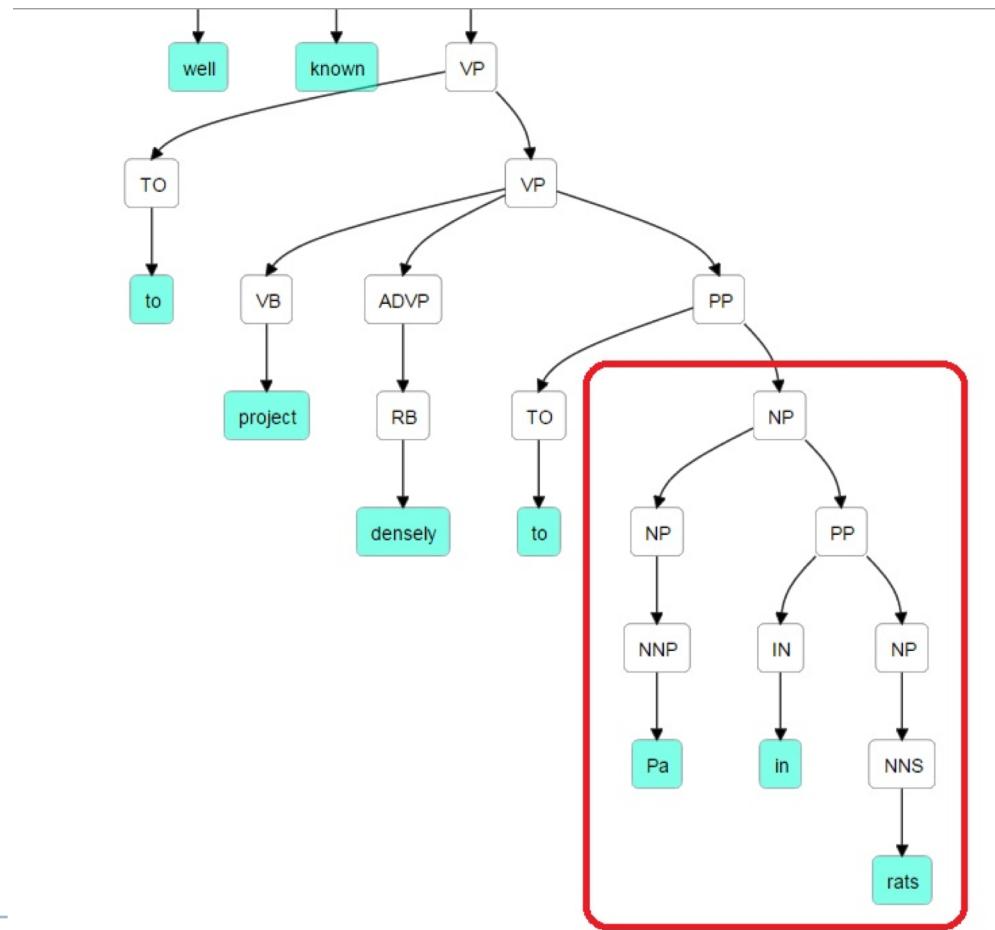
# İlişki tespiti

---

- Sözdizim ve bağlam ağaçları kullanarak cümlelerin gramer analizi yapıldı
  - İlişkinin yönünün tespiti
  - Etkileşen ve etkilenenlerin tespiti

# Sözdizim Ağacı

*“The suprachiasmatic nucleus is well known to project densely to Pa in rats”*



# Örnek uygulama: PVT

- Elle etiketlenen PVT derlemi

- 14 tam metinli makale
- 322 ilişki

- PubMed'deki tüm PVT ile ilgili makaleler

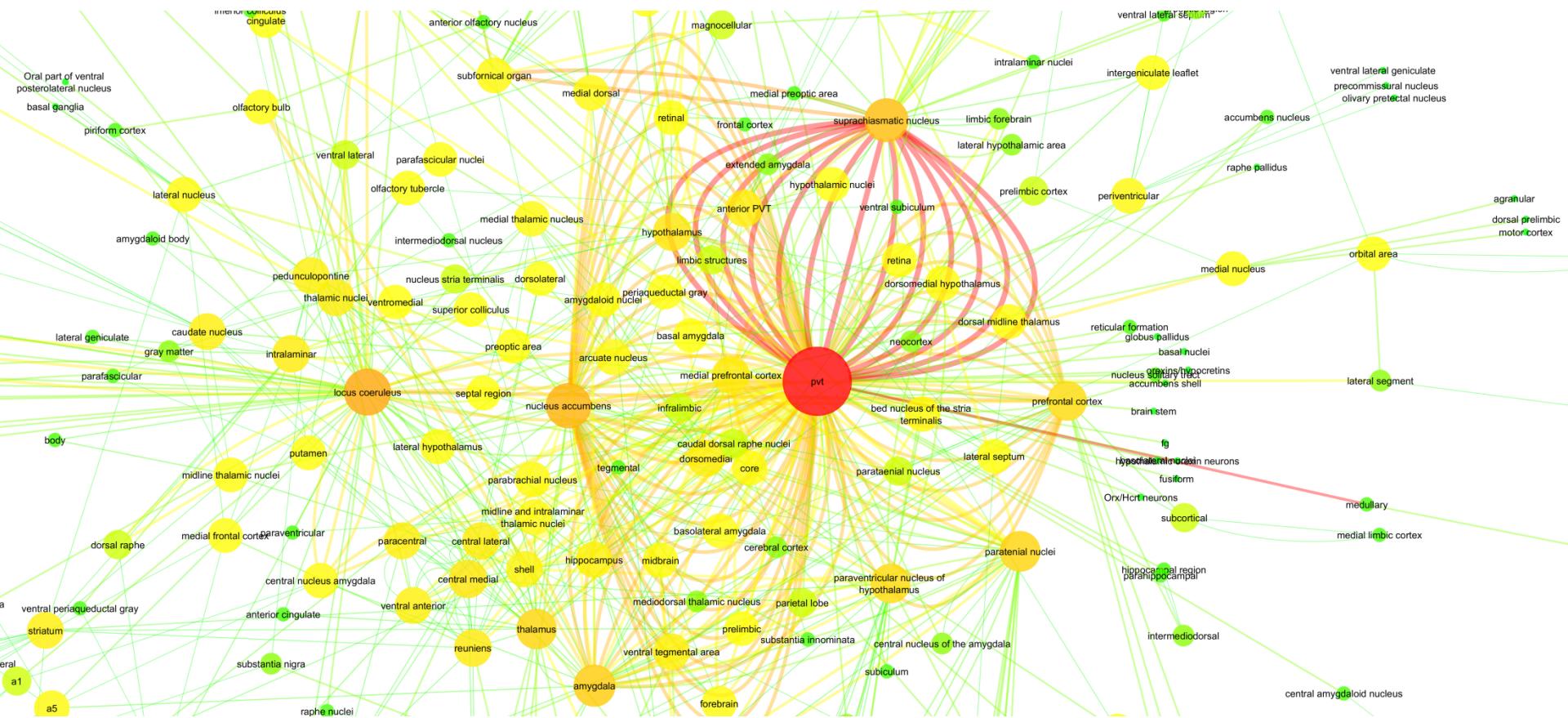
- 451 özet
- 107 tam metinli makale

# PVT Derlemi Üzerinde Değerlendirme

---

	Precision	Recall	F-Measure
Strict (Full Match)	66.43%	33.23%	44.30%
Lenient (Full Match + Partial Match)	75.78%	37.89%	50.52%
NLP-based	87.58%	43.79%	58.39%

# PVT Bağlantı Ağacı



# Yeni hipotez

---

- ▶ PVT aşağıdaki beyin bölgeleri ile kuvvetli bağlantılar kuruyor
  - ▶ SCN, nucleus accumbens, amygdaloid complex
  - ▶ extended amygdala incl.
    - ▶ bed nucleus of the stria terminalis
    - ▶ ventromedial prefrontal cortex
- ▶ Bu beyin kısımları depresyon ve ruh hali ile ilişkilendirilmiştir.
- ▶ PVT'nin de depresyon mekanizmasında önemli bir rolü olabilir.
  - ▶ Only, Zhu et al. (2011) suggests that PVT neurons might be engaged in acute depressive events

---

# Varlık İsmi Normalizasyonu

İşbirliği: İlknur Karadeniz

Karadeniz, İ., & Özgür, A. (2015). **Detection and categorization of bacteria habitats using shallow linguistic analysis.** *BMC bioinformatics*, 16(10), S5.

Karadeniz, İ., & Özgür, A. (2019). **Linking entities through an ontology using word embeddings and syntactic re-ranking.** *BMC bioinformatics*, 20(1), 156.

---



# Varlık İsmi Normalizasyonu

id:OBT:000164  
name: **respiratory tract**  
synonym:  
"respiratory tree"  
synonym:  
"respiratory"  
is\_a: OBT:000065 !  
animal part

The etiologic and epidemiologic spectrum of bronchiolitis in **pediatric** practice.

To develop a broad understanding of the causes and patterns of occurrence of wheezing associated **respiratory** infections, we analyzed data from an 11-year study of acute lower **respiratory** illness in a **pediatric** practice. Although half of the WARI occurred in **children less than 2 years of age**, wheezing continued to be observed in 19% of **children greater than 9 years of age** **who had lower respiratory illness**. **Males** experienced LRI 1.25 times more often than did **females**; the relative risk of **males** for WARI was 1.35. A nonbacterial pathogen was recovered from 21% of **patients with WARI**; respiratory syncytial virus, parainfluenza virus types 1 and 3, adenoviruses, and Mycoplasma pneumoniae accounted for 81% of the isolates. **Patient** age influenced the pattern of recovery of these agents. The most common cause of WARI in **children under 5 years of age** was RSV whereas Mycoplasma pneumoniae was the most frequent isolate from **school age children with wheezing illness**. The data expand our understanding of the causes of WARI and are useful to **diagnosticians** and to **researchers** interested in the control of lower **respiratory** disease.

id:OBT:002307  
name: **pediatric patient**  
is\_a: OBT:002133 ! Patient  
is\_a: OBT:002146 ! child

# Zorluklar

---

The etiologic and epidemiologic spectrum of bronchiolitis in **pediatric** practice.

To develop a broad understanding of the causes and patterns of occurrence of wheezing associated **respiratory** infections, we analyzed data from an 11-year study of acute lower **respiratory** illness in a **pediatric** practice. Although half of the WARI occurred in children less than 2 years of age, wheezing continued to be observed in 19% of **children greater than 9 years of age** who had lower respiratory illness. Males experienced LRI 1.25 times more often than did females; the relative risk of males for WARI was 1.35. A nonbacterial pathogen was recovered from 21% of **patients with WARI**; respiratory syncytial virus, parainfluenza virus types 1 and 3, adenoviruses, and Mycoplasma pneumoniae accounted for 81% of the isolates. **Patient** age influenced the pattern of recovery of these agents. The most common cause of WARI in **children under 5 years of age** was RSV whereas Mycoplasma pneumoniae was the most frequent isolate from **school age children with wheezing illness**. The data expand our understanding of the causes of WARI and are useful to **diagnosticicians** and to **researchers** interested in the control of lower **respiratory** disease.

Sözcüksel  
benzerlik yok!

id:OBT:002307  
name: **pediatric patient**  
is\_a: OBT:002133  
! Patient  
is\_a: OBT:002146  
! child

# Zorluklar – Belirsizlik

Doğru kavram

id:OBT:002307  
name: **pediatric patient**

is\_a: OBT:002133  
! Patient  
is\_a: OBT:002146  
! child

id: OBT:002167  
name: **boy**  
is\_a: OBT:002146  
! child

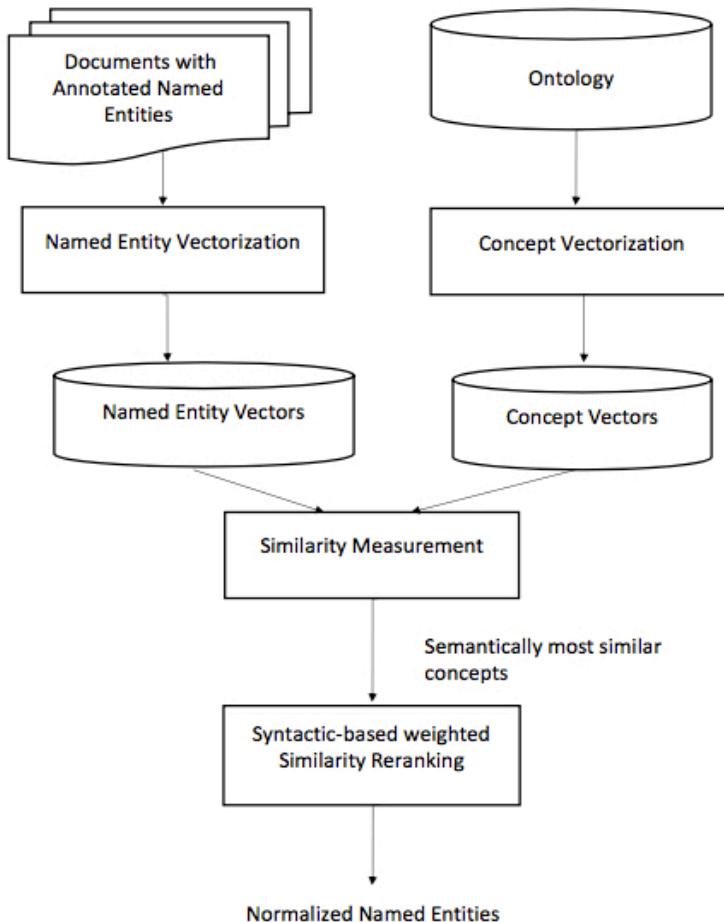
The etiologic and epidemiologic spectrum of bronchiolitis in **pediatric** practice.

To develop a broad understanding of the causes and patterns of occurrence of wheezing associated **respiratory** infections, we analyzed data from an 11-year study of acute lower **respiratory** illness in a **pediatric** practice. Although half of the WARI occurred in **children less than 2 years of age**, wheezing continued to be observed in 19% of **children greater than 9 years of age who had lower respiratory illness**. Males experienced LRI 1.25 times more often than did **females**; the relative risk of **males** for WARI was 1.35. A nonbacterial pathogen was recovered from 21% of **patients with WARI**; respiratory syncytial virus, parainfluenza virus types 1 and 3, adenoviruses, and *Mycoplasma pneumoniae* accounted for 81% of the isolates. **Patient** age influenced the pattern of recovery of these agents. The most common cause of WARI in **children under 5 years of age** was RSV whereas *Mycoplasma pneumoniae* was the most frequent isolate from **school age children with wheezing illness**. The data expand our understanding of the causes of WARI and are useful to **diagnosticians** and to **researchers** interested in the control of lower **respiratory** disease.

Sözcüksel olarak en çok benzeyen kavram.

id: OBT:000608  
name: **male animal**  
synonym: "male"  
is\_a: OBT:000380  
! animal with age or sex property

# Sistemin Genel Yapısı



Denetimsiz bir yaklaşım

Elle etiketli veri veya alana özel sözlüklerde ihtiyaç duymuyor.

Başka alanlara uyarlanabilir.

# Varlık isimlerinin vektörel gösterimi

day

$$\vec{e}(\text{day}) = [-0.25875682 \ -0.11159618 \dots 0.360897]$$

care

$$\vec{e}(\text{care}) = [0.18755111 \ 0.35023546 \dots 0.3920108]$$

center

$$\vec{e}(\text{center}) = [0.20950332 \ 0.5162147 \dots -0.13443379]$$

$$\overrightarrow{\text{sum}} = \vec{e}(\text{day}) + \vec{e}(\text{care}) + \vec{e}(\text{center}) =$$

$$[0.13829761 \ 0.75485398 \dots 0.61847401]$$

Divide by the  
number of  
tokens n in the  
phrase  
(for the sample  
 $n = 3$ )

$$\vec{e}(\text{a day-care center}) = \overrightarrow{\text{sum}} / 3 =$$
$$(0.13829761/3) \ (0.75485398/3) \dots (0.61847401/3)]$$

a day-care  
center

$$\vec{e}(\text{a day-care center}) =$$
$$[0.0460992 \ 0.25161799 \dots 0.206158]$$

# Örnek bir yanlış pozitif

---

Table 1. Semantically most similar concepts to the entity mention "children attending a day-care center" without reranking

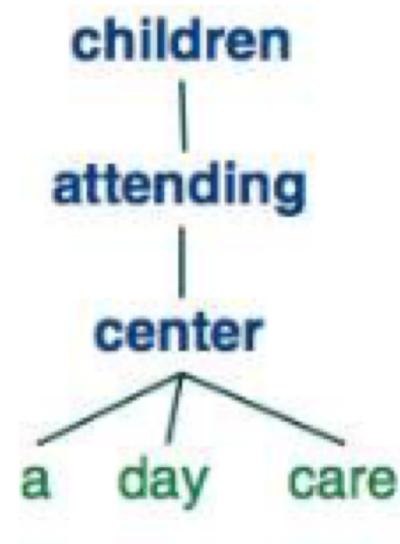
Ranking	Concept	Similarity score
1	OBT:001423 medical center	0.8297
2	OBT:001801 clinic	0.7917
28	OBT:002146 child	0.6979

---

# Yeniden sıralama

---

- ▶ Cümlelerin sözdizimsel analizi kullanıldı
- ▶ Aday varlık isminin baş kelimesinin (headword) bulunması
- ▶ Aday kavramın baş kelimesinin bulunması
- ▶ Baş kelimeler arasındaki benzerlige ağırlık verilmesi



$$S_{RR}(m, c) = (w * S_S(m_{head}, c_{head})) + ((1-w)*S_S(m, c))$$

# Yeniden sıralamanın etkisi

Table 1. Semantically most similar concepts to the entity mention "children attending a day-care center" without reranking

Ranking	Concept	Similarity score
1	OBT:001423 medical center	0.8297
2	OBT:001801 clinic	0.7917
28	OBT:002146 child	0.6979

Table 2. Semantically most similar concepts to the entity mention "children attending a day-care center" after syntactic reranking

Ranking	Concept	Similarity score
1	OBT:002146 child	0.7484
3	OBT:001801 clinic	0.6519
24	OBT:001423 medical center	0.5460

System	Train	Dev
<b>Before Re-ranking</b>	0.601	0.629
<b>After Re-ranking</b>	0.648	0.677

# Sonuçlar

---

System	Precision
<b>BOUNEL</b>	<b>0.659</b>
<b>TURKU</b>	0.630
<b>BOUN</b>	0.620
<b>CONTES</b>	0.597
<b>LIMSI</b>	0.438

---

# Biyomedikal alanda cümleler arası semantik benzerlik hesaplama

İşbirliği: Gizem Soğancıoğlu ve Hakime Öztürk

G. Soğancıoğlu, H. Öztürk, A. Özgür. "**BIOSSES: A Semantic Sentence Similarity Estimation System for the Biomedical Domain**", Bioinformatics, 2017.

---



# Biyomedikal alandan örnek

## Semantically Highly Similar Sentence Pairs

- S1: This form of necrosis, also termed necroptosis, requires the activity of receptor-interacting protein kinase 1 and its related kinase 3 .
- S2: Moreover, other reports have also shown that necroptosis could be induced via modulating RIP1 and RIP3.

Amaç:

Semantik cümle benzerliği hesaplamak için bir sistemin geliştirilmesi.

# Benzerlik Yöntemleri

---

- ▶ Karakter dizisi benzerliği
- ▶ Dağıtık temsil vektörü benzerliği
- ▶ Ontoloji tabanlı benzerlik
- ▶ Denetimli makine öğrenmesiyle farklı benzerlik ölçütlerinin birleştirilmesi



# Karakter dizisi benzerliği



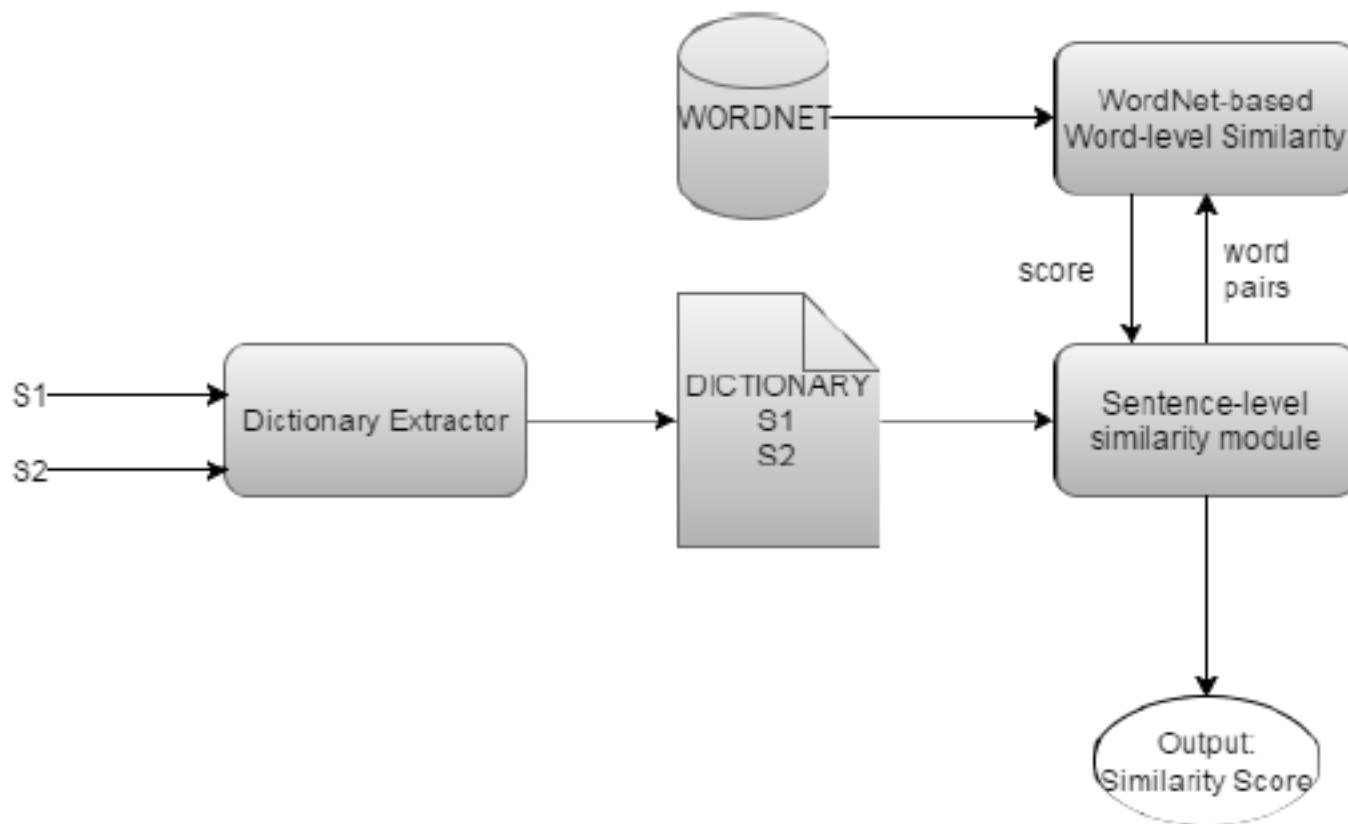
# Ontoloji tabanlı yöntemler

---

- ▶ WordNet
- ▶ UMLS
- ▶ Birleştirilmiş Benzerlik Ölçütü



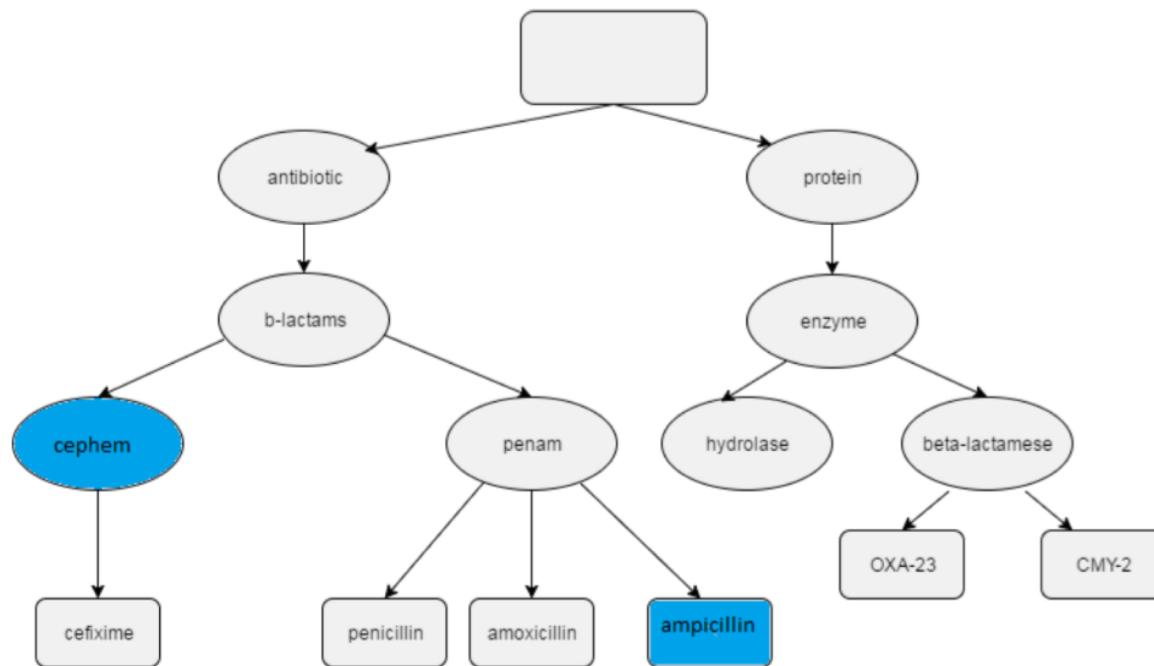
# Ontoloji-tabanlı benzerlik



# Örnek: Ontoloji-tabanlı benzerlik

$$\text{Sim\_Path}(c_1, c_2) = 2 * \maxDepth - \text{len}(c_1, c_2)$$

$$\text{Sim\_Path}(\text{cephem}, \text{ampicillin}) = 10 - 4 = 6$$



## Örnek

---

$S_1$  = Necroptosis requires the activity of RIP1 and RIP3.

$D_1$  = [1,1,1,1,1,1,1,0,0,0.33,0,0.1]

$S_2$  = Necroptosis could be induced via modulating RIP1 and RIP3.

$D_2$  = [1,0.33,0,0.1,0,1,1,1,1,1,1,1]

Standart:

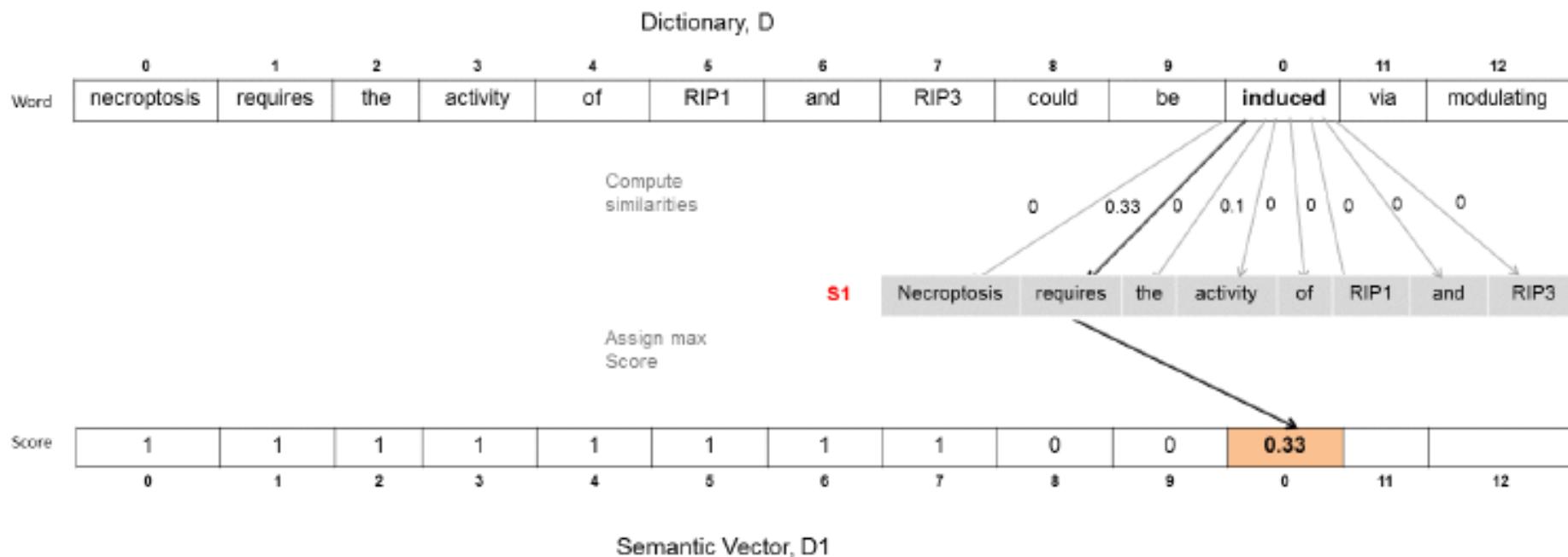
Kosinüs benzerliği( $S_1, S_2$ )= 0.47

Ontoloji-tabanlı:

Kosinüs benzerliği( $D_1, D_2$ ) = 0.56

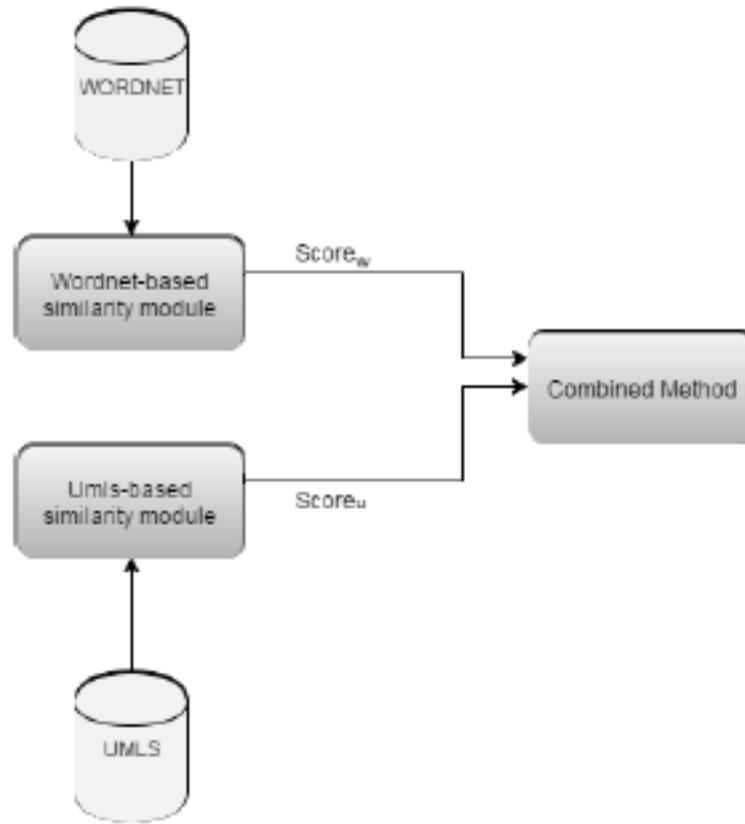


# Örnek

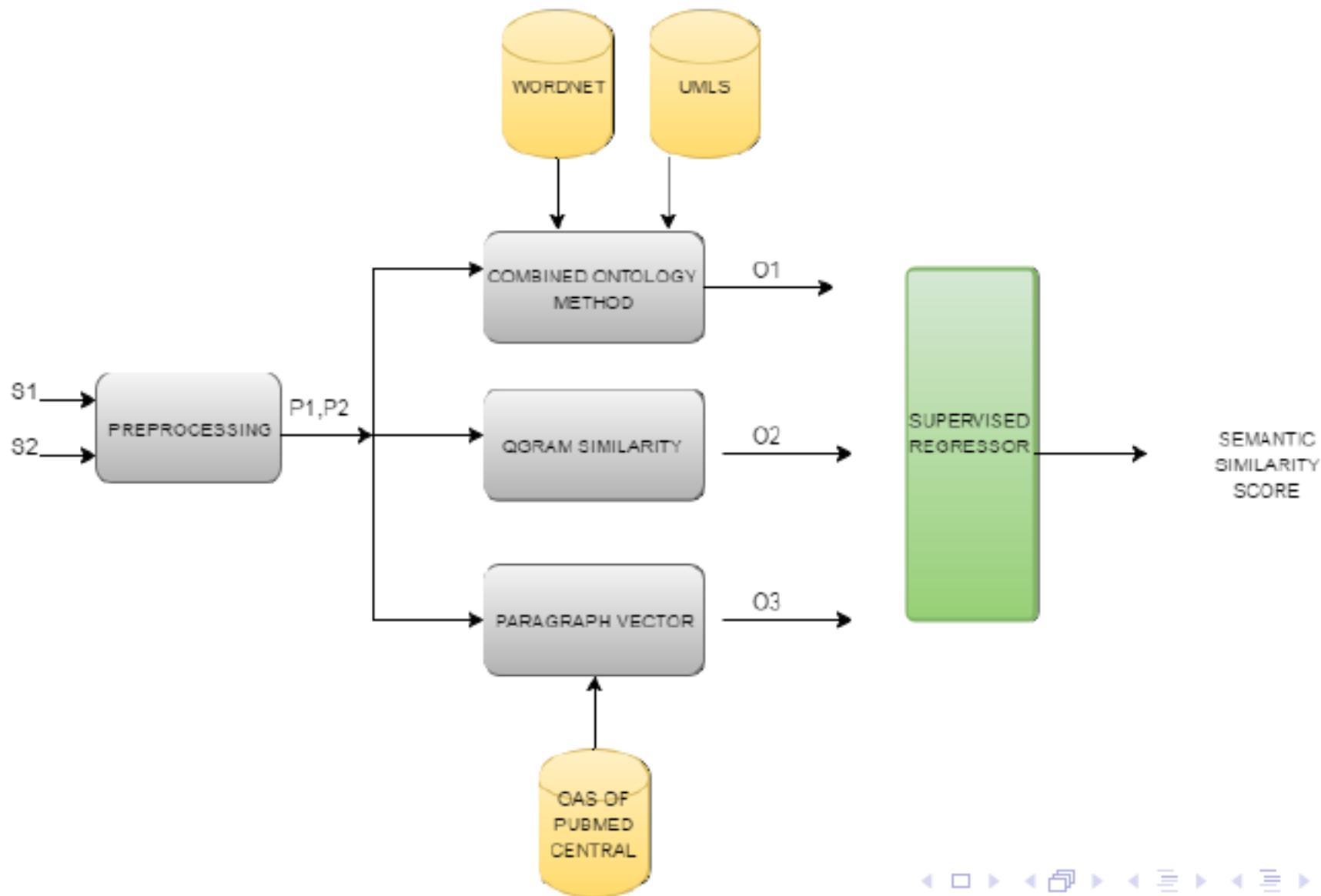


# Birleştirilmiş skor

$$\text{combined score} = \text{Score}_{\text{WordNet}} \cdot \lambda + \text{Score}_{\text{UMLS}} \cdot (1 - \lambda)$$



# Denetimli Yaklaşım



# Veri kümesi etiketlenmesi

---

<b>Annotation Score</b>	<b>Definition</b>
<b>0</b>	The two sentences are on different topics.
<b>1</b>	The two sentences are not equivalent, but are on the same topic.
<b>2</b>	The two sentences are not equivalent, but share some details.
<b>3</b>	The two sentences are roughly equivalent, but some important information differs/missing.
<b>4</b>	The two sentences are completely or mostly equivalent, as they mean the same thing.



# Cümle seçimi

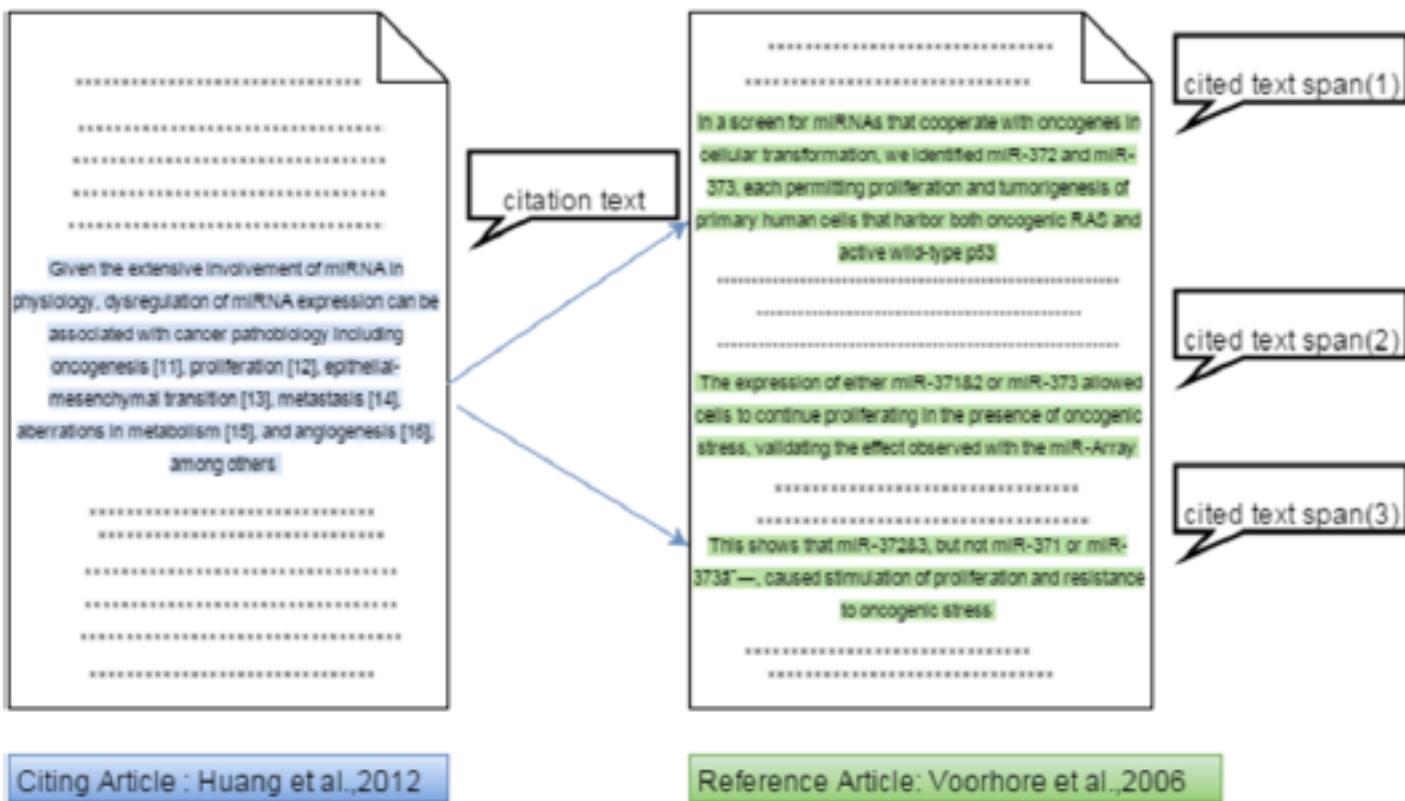


Figure: Example Annotation from TAC<sup>2</sup> training data set

Pair\_ID|Sentence1|Sentence2|Similarity\_Score|Annotator\_ID

# Sonuçlar

---

Table: Correlation scores among annotators

	Correlation r
Annotator A	0.952
Annotator B	0.958
Annotator C	0.917
Annotator D	0.902
Annotator E	0.941

Table: Correlation scores for domain-independent state-of-the-art systems

	Correlation r
ADW	0.586
SEMIAR	0.419



# Sonuçlar

Methods	Pearson Correlation
<b>Domain-independent Systems</b>	
ADW	<b>0.586</b>
SEMILAR	0.419
<b>String Similarity Measures</b>	
Qgram	<b>0.754</b>
<b>Word Embeddings based Similarity</b>	
Paragraph Vector	<b>0.787</b>
<b>Ontology-based Similarity</b>	
WordNet based Similarity Module-Path	<b>0.644</b>
UMLS based Similarity Module-Path	<b>0.651</b>
Combined Ontology Method( $[\lambda = 0.5]$ )	<b>0.710</b>
Supervised Model - Linear Regression	<b>0.836</b>





# BIOSSES : Biomedical Semantic Similarity Estimation System

Method:

Qgram Similarity (Lexical)

Necroptosis requires the activity of RIP1 and RIP3.

Necroptosis could be induced via modulating RIP1 and RIP3.

Calculate

0.46846848726272583

About    Data Set    Source Code

This web site is best viewed in Google Chrome and Internet Explorer.

BioSSES computes similarity of biomedical sentences by utilizing WordNet as the general domain ontology and UMLS as the biomedical domain specific ontology.

We allow you to compute sentence similarity with the following methods:

- Qgram [0-1]
- Wordnet [0-1]
- UMLS [0-1]
- Paragraph Vector [0-1]
- Combined Ontology (Wordnet and UMLS) [0-1]
- Supervised Approach [0-4]

For citing this study, please use:

BioSSES is a set of Java codes for computing semantic sentence similarity, developed by [Gizem Sogancioglu](#), [Hakime Ozturk](#) and [Arzucan Ozgur](#), in the Department of [Computer Engineering](#), Bogazici University.

BioSSES is open-source software made available under the terms of the [The GNU Common Public License v.3.0](#). You are free to use the code under those terms.

# Teşekkür

---

## Öğrencilerim

İlknur Karadeniz (İşık Ü. Öğretim Üyesi)  
Hakime Öztürk  
N. Özlem Özcan Şimşek  
Erinç Gökdeniz  
Ferhat Aydin  
Zehra Melce Hüsünbeyi  
Gizem Soğancıoğlu

## İşbirliği Yaptığımız Araştırmacılar

Elif Özkırımlı (Boğaziçi Ü.)  
Reşit Canbeyli (Boğaziçi Ü.)  
Fikret Gürgen (Boğaziçi Ü.)  
Kutlu Ülgen (Boğaziçi Ü.)  
Yongqun Oliver He (University of Michigan)  
Junguk Hur (University of North Dakota)  
Dragomir Radev (Yale University)

## Destekler

EU Marie Curie Career Integration Grant  
Bilim Akademisi Genç Bilim İnsanları Ödül Programı (BAGEP)  
TUBITAK-BIDEB 2211  
Boğaziçi Ü. BAP (No 12304)  
Boğaziçi Ü. BAP (No D-13242)  
DPT TAM Projesi (No 2007K120610)



# Teşekkürler!

---

