

# Machine-Learning for chemogenomics: Application to target identification for Cystic Fibrosis

**Or : How to identify the targets of drugs identified in phenotypic screens ?**

Pr Isabelle Sermet (Institut Necker Enfants Malades, INEM)

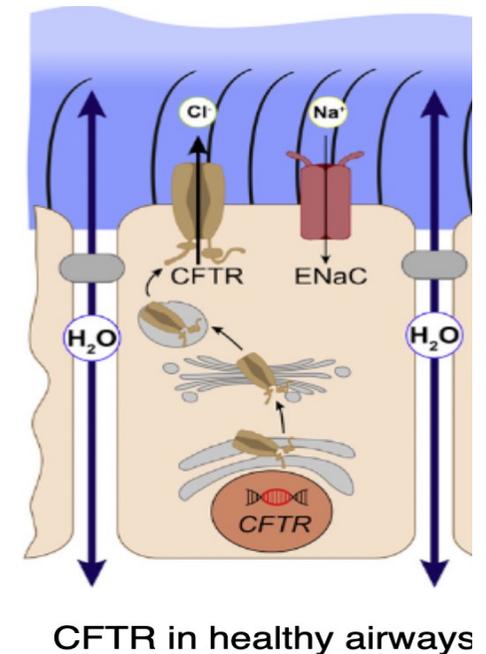
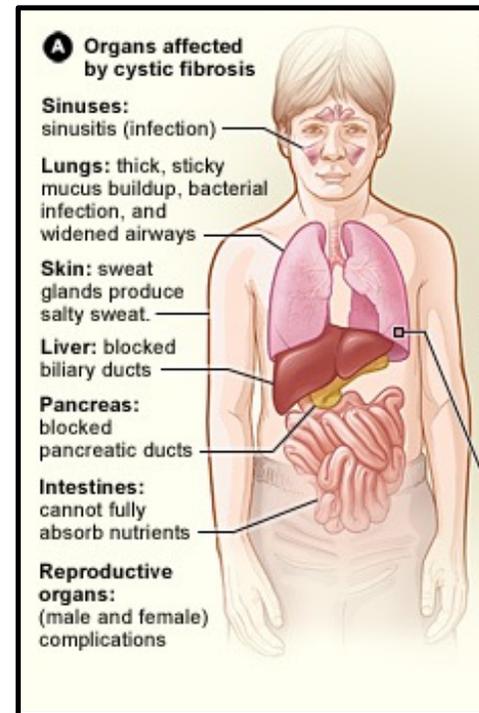
Pr Véronique Stoven (CBIO, Mines Paris PSL)

Laurence Calzone, Loredana Martignetti (U1331, Institut Curie)



# Cystic Fibrosis (CF)

- Most frequent genetic disease in caucasian population (1/3500 birth)
- **Mutation in *cftr* gene** (Cystic Fibrosis Transmembrane Regulator) coding for **CFTR protein**: Cl<sup>-</sup> channel at the apical surface of epithelial cells
- **> 2000 mutations**, various **impacts** on CFTR protein (classes)
- **70% cases: deletion of F508**, folding defect in CFTR, elimination via proteasome, absence at membrane (Class II)
- **Pulmonary symptoms**: chronic infections, inflammation, tissue damage, altered pulmonary function



**Until recently** (10 years): symptomatic treatments (antibiotics, physiotherapy, nutritional support etc...)

## CFTR modulators : proteic therapy

**CFTR Activators and Correctors** : improve processing, maturation and channel activity of mutated CFTR

### Limitations :

- Patients bearing **at least one F508del mutation**
- Some patients are **non responders**
- **15% patients** bear “un-rescuable” mutations → **not eligible**.
- **Mechanism of action not fully understood**: identified in phenotypic screens (increase of chloride conductance).

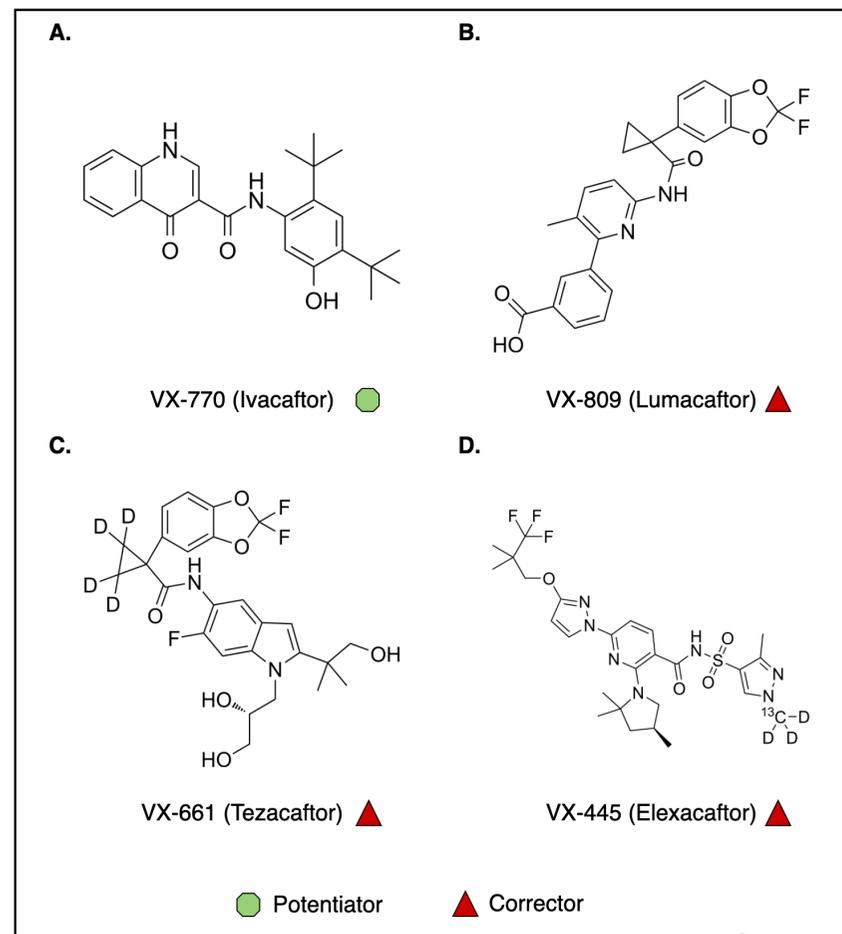
### Issues:

**Identify ETI protein targets** other than CFTR

(clinical response not correlated with restoration of Cl<sup>-</sup> conductance)

**Propose new therapeutic targets** and their associated “mutation agnostic” drugs (ETI or optimized for them).

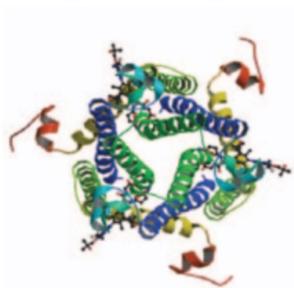
**Trikafta™** (ETI, 4 years)= 2 correctors (Elexacaftor, Tezacaftor) + 1 activator (Ivacaftor)



## Identify new therapeutic targets among ETI binding proteins

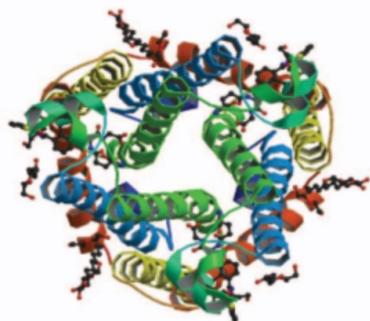
### Deleterious off-target (side-effects)

OFF Target Protein  
(side-effect)



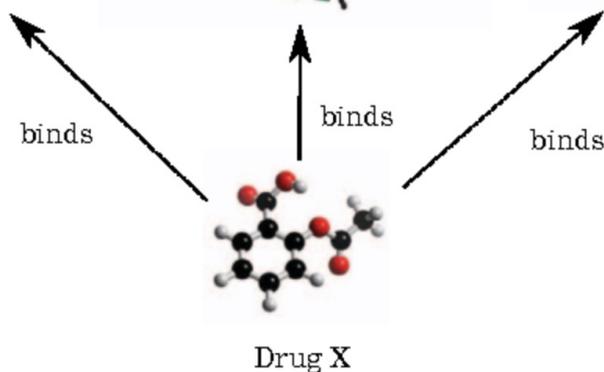
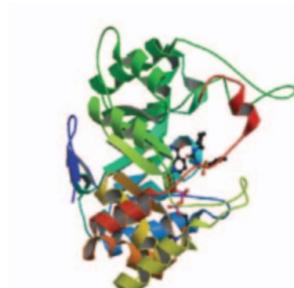
### CFTR

ON Target Protein  
(Target Disease)



### Beneficial off-target = therapeutic targets

OFF Target Protein  
(side-effect)



**Difficult problem** : 25.000 human genes, 3000 known druggable proteins

→ Predict ETI protein interaction profile

**Boils down to a chemogenomic problem:**

Predict drug-target interactions (DTI) at large scale

**Of general interest in drug discovery**

Phenotypic drugs candidates in cancer (survival assays)

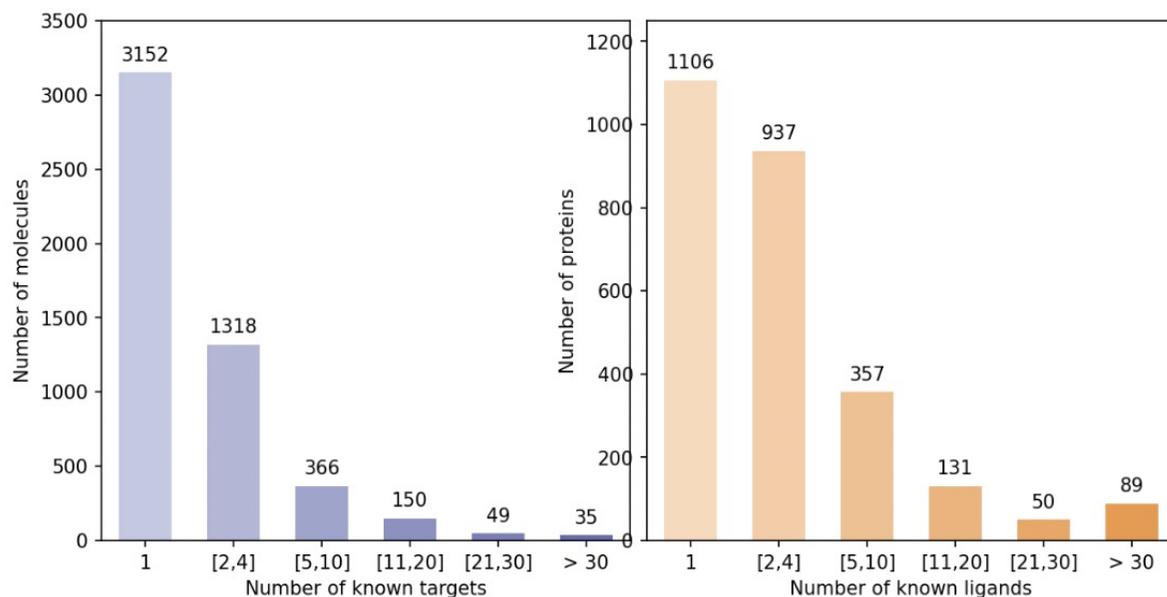
Anticipate side-effects, drug repurposing etc...

Using a **ML algorithm** and a **training set of DTIs**

## Quality of training set (known interactions): critical for DTI prediction with ML

### Limitations in classical DTI training sets:

- bias for well-studied families of proteins (kinases, GPCRs, NRs etc...). Idem for mol
- no negative interactions. **Randomly chosen among unknown interactions**



**DrugBank : 8.600 interactions out of 12.210 involve proteins with > 10 known ligands**

Protein nb of Ligands	nb of Interactions
1	1106
2 to 4	2527
5 to 10	2404
11 to 20	1920
21 to 30	1238
>30	5442

**Proteins with many ligands viewed by ML algorithm as probable targets:  
leads to False Positive predictions**

Nb of known ligands	FPR
Prot in Category	RN-Datasets
0	2.2 ± 0.4
1	3.7 ± 0.5
2 to 4	5.1 ± 0.9
5 to 10	9.9 ± 0.9
11 to 20	13.8 ± 1.7
21 to 30	23.0 ± 4.9
> 30	18.6 ± 2.8

## Correct bias by “balanced” positive and negative DTIs per protein and molecule

Nb of known ligands

Prot in Category	FPR (Threshold = 0.5 )	
	RN-Datasets	BN-Datasets
0	2.2 ± 0.4	3.1 ± 0.5
1	3.7 ± 0.5	3.1 ± 0.8
2 to 4	5.1 ± 0.9	6.4 ± 0.8
5 to 10	9.9 ± 0.9	8.3 ± 0.6
11 to 20	13.8 ± 1.7	10.6 ± 0.5
21 to 30	23.0 ± 4.9	12.0 ± 3.0
> 30	18.6 ± 2.8	9.0 ± 0.4

Performance of ML algorithm according to False Positive Rate (FDR) on DrugBank

randomly choose negatives such that **equal nb of positives and negatives for each protein and molecule**

**Target prediction improvement:**  
**Improves score and rank of true target for 3 marketed drugs**

Drug	RN-Datasets		BN-Datasets	
	Target Score	Target Rank	Target Score	Target Rank
DB11363	0.8	31	0.8	3
DB11842	0.76	31	0.85	18
DB11732	0.67	107	0.83	17

# Build a large high-quality DTI training dataset : LCldb

## Binary interaction databases

**DrugBank**: DTI between FDA approved or in development molecules and their targets

Drugbank v1.5.1 [Wishart *et al*, 2018]

2.513 proteins

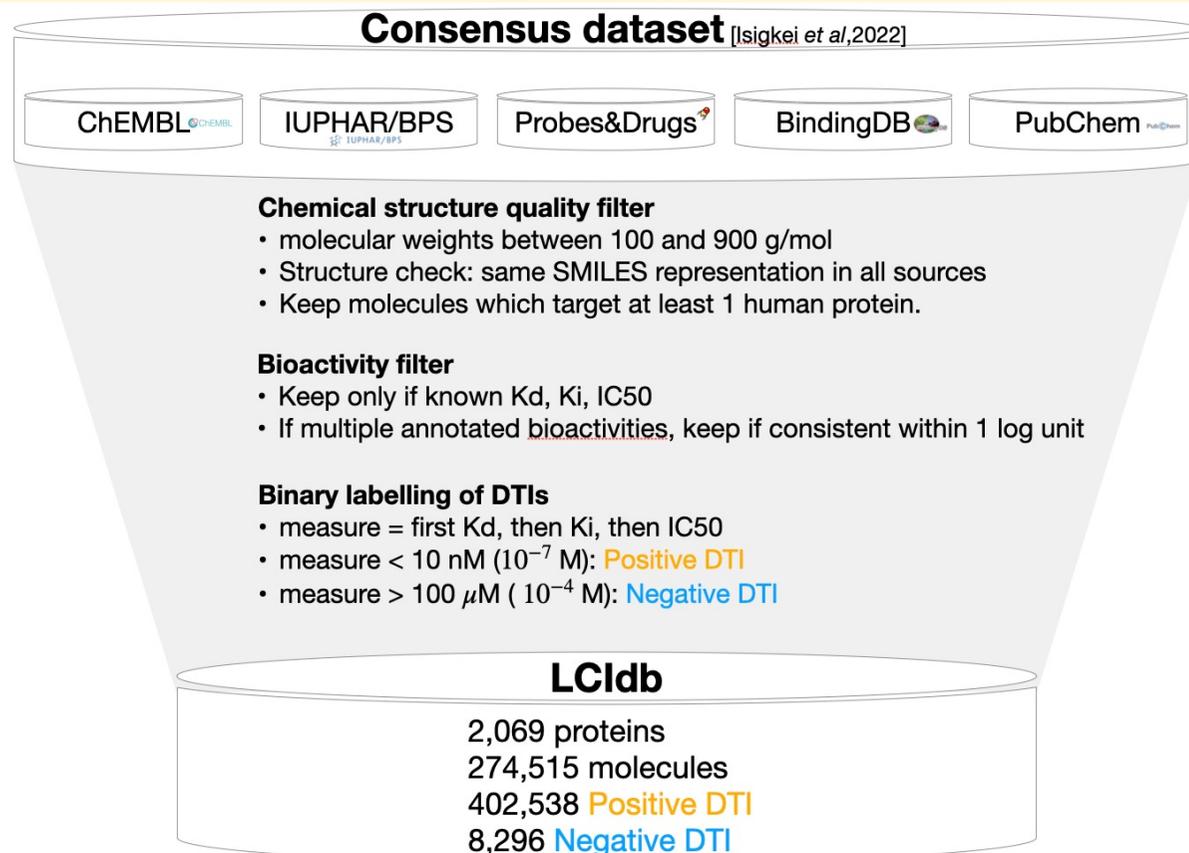
4.813 molecules

13.716 + interactions ■

BIOSNAP similar size

BindingDB smaller

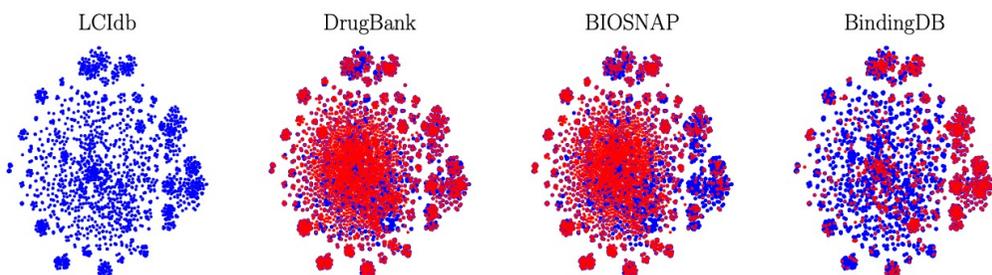
## Bioactivity databases



Complete with random “balanced” negatives per protein and molecule

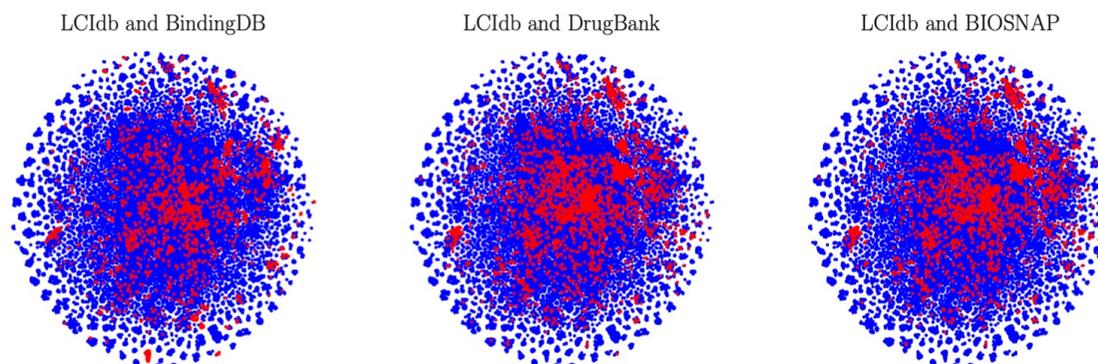
## Characterisation of LCIdb

### 2D representation of protein space coverage (t-SNE algorithm on protein features)



Datasets	Molecules	Proteins	Positive DTIs
BIOSNAP	<u>4,510</u>	<u>2,181</u>	<u>13,836</u>
Unseen_drugs			13,836
Unseen_targets			13,836
BindingDB	<u>7,161</u>	<u>1,254</u>	<u>9,166</u>
DrugBank	<u>4,813</u>	<u>2,507</u>	<u>13,715</u>
DrugBank (Ext)	<u>4,257</u>	<u>1,216</u>	<u>10,838</u>
LCIdb	<u>271,180</u>	<u>2,060</u>	<u>396,798</u>
Unseen_drugs	271,180	2,060	396,798
Unseen_targets	271,180	2,060	396,798
Orphan	191,901	2,060	208,041

### 2D representation of chemical space coverage (t-SNE algorithm on molecule features)



### LCIdb

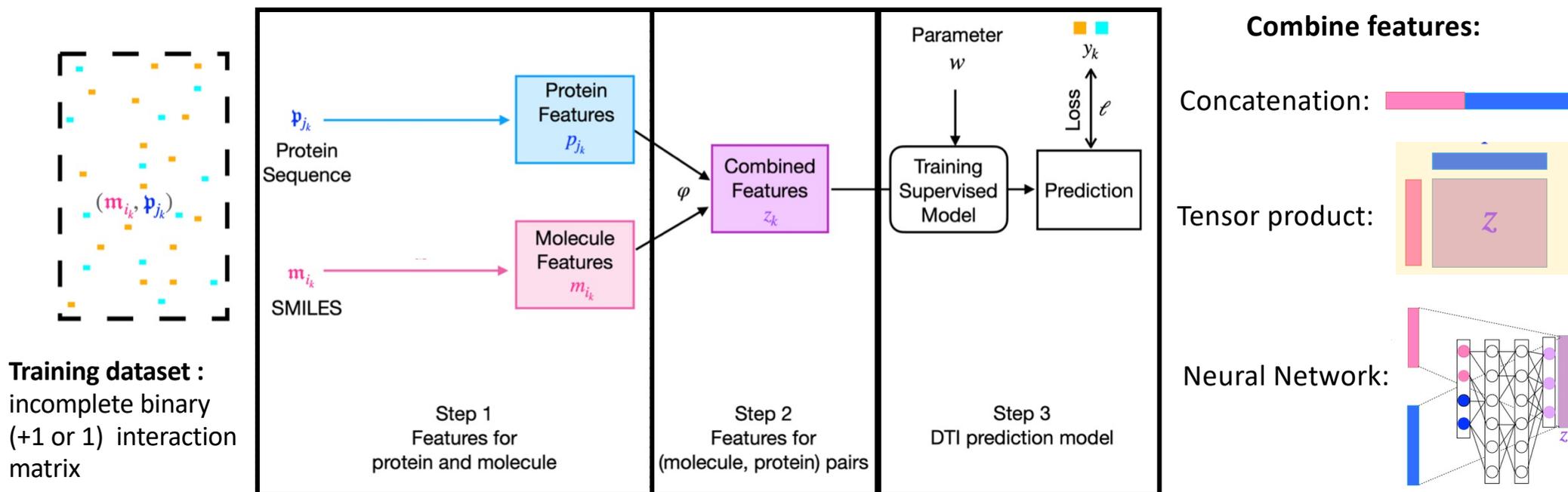
**30 times larger than BIOSNAP and DrugBank**  
**43 times larger than BindingDB**

Similar nb of proteins (“druggable” human proteins)  
Much larger nb of molecules

*Guichaoua et al, JCIM 2024*

## General pipeline for DTI prediction in ML chemogenomic framework

A classification problem: distinguish **interacting** and **not interacting**  $(m, p)$  pairs



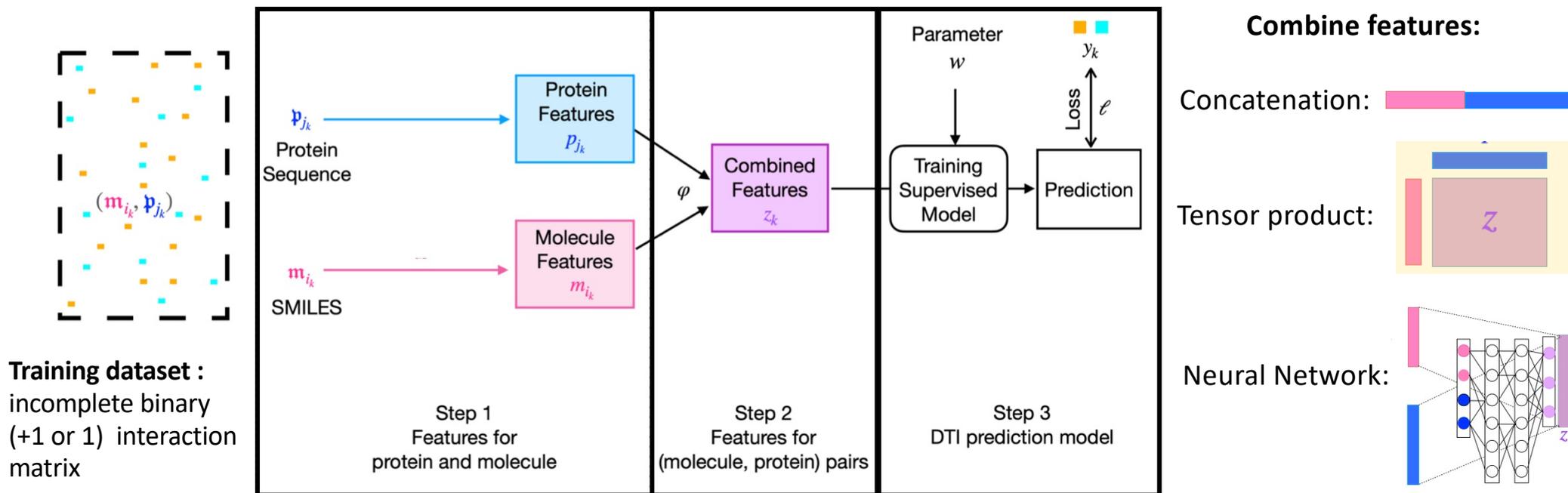
Computational burden for large training sets (LCIdb)

### KOMET : efficient ML algorithm

- Step 1: molecule and protein **features derived from Kernel methods**
- Steps 2-3: **tensor product** with “classical” SVM algorithm in the space of  $(m, p)$  pairs

## General pipeline for DTI prediction in ML chemogenomic framework

A classification problem: distinguish **interacting** and **not interacting**  $(m, p)$  pairs



Computational burden for large training sets (LCIdb)

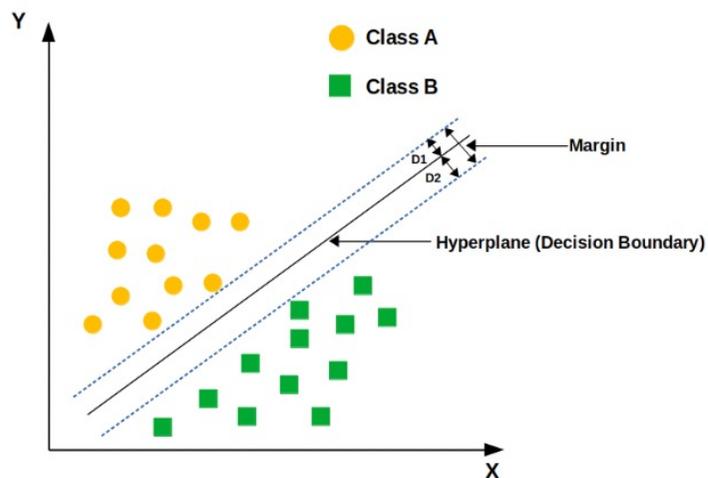
### KOMET : efficient ML algorithm

- Step 1: molecule and protein **features derived from Kernel methods**
- Steps 2-3: **tensor product** with “classical” SVM algorithm in the space of  $(m, p)$  pairs

# Train a SVM = Learn hyperplane separating positive/negative data points

Training set:  $(x_i, y_i)$  with  $y_i = \mp 1$  (labels)

## Linearly separable case

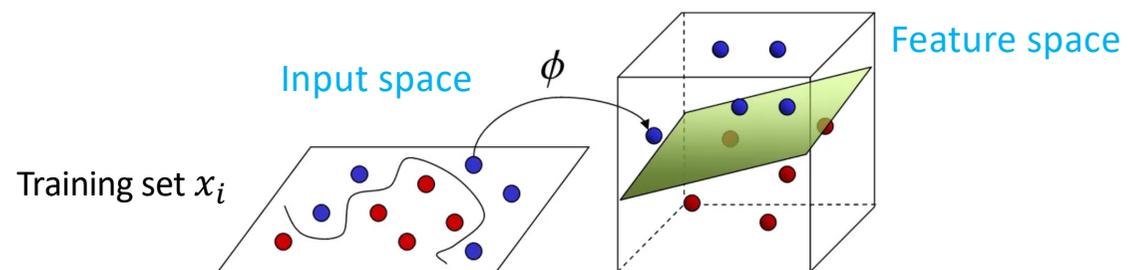


compute hyperplane with largest margin :  
only involves  $\langle x_i, x_j \rangle$

→ **Good generalization properties**

**Prediction based on position of points**

**Linearly non-separable case:** may be linearly separable after transformation  $\phi$  in a feature space (of higher dimension)



**Kernel trick** provides a solution to this problem:

- **define  $k$**  : a similarity measure (kernel) in the input space
- **corresponds to a mapping function  $\phi$**  in a feature space such that:  
 $\langle \phi(x_i), \phi(x_j) \rangle = k(x_i, x_j) \longrightarrow$  kernel representation of data
- compute hyperplane only based on  $k(x_i, x_j)$  ( **$\phi$  not computed explicitly**)

**Predictions based on  $k(x, x_i)$**

## The “classical” approach for SVM with kernel methods for chemogenomics

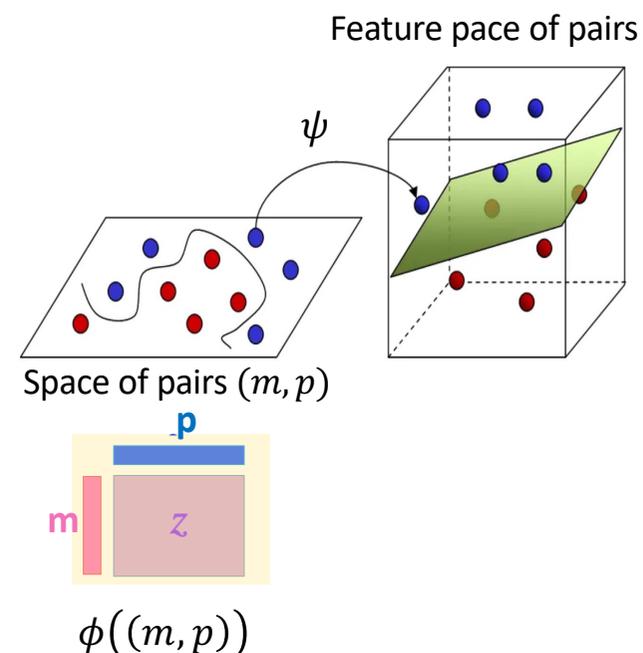
Training set :  $(m, p)$  pairs (with labels +1 or -1)

- choose encodings for molecules and proteins  $\phi_M(m)$  and  $\phi_P(p)$  (i.e. descriptors)
- define similarity measures (kernels) on molecules and proteins :  $k_M$  and  $k_P$
- choose tensor product to encode  $(m, p)$  pairs:  $\phi((m, p)) = \phi_M(m) \otimes \phi_P(p)$
- hyperplane in the feature space of pairs computed based only on  $k_M$  and  $k_P$   
( based on  $k_M(m_i, m_j) \cdot k_P(p_i, p_j)$ ,  $\phi((m, p))$  not computed)

Motivation to use Tensor product with SVM kernel methods :

- Systematic definition and good “mixing properties”
- Efficient computation

Predictions only based on  $k_M(m, m_i) \cdot k_P(p, p_i)$



**Limitation: computational burden for training on large DTIs datasets (LCIdb)**

## Principle of the KOMET algorithm: conceived for large training sets

KOMET solves a “classical SVM optimization problem” in the feature space of  $(m, p)$  pairs based on:

- “kernel-derived” encoding of molecules and proteins  $\phi_M$  and  $\phi_P$
- $(m, p)$  pairs encoded as tensor product:  $\phi((m, p)) = \phi_M(m) \otimes \phi_P(p)$
- “classical” SVM optimisation in  $(m, p)$  space

Compute protein mapping function  $\phi_P$  :

- choose a similarity measure  $k_P$  defining matrix  $K_P$  as:  $(K_P)_{ij} = k_P(p_i, p_j)$
- define  $X_P$  such as:  $K_P = X_P \cdot X_P^T$
- compute matrix  $X_P$  : defines kernel-derived encoding  $\phi_P$

$X_P$  is computed by SVD decomposition of  $K_P$

Computation of  $K_P$  and  $X_P$  “feasible” : size  $n_P \cdot n_P$  ( $n_P = 2060$  in LCIdb)

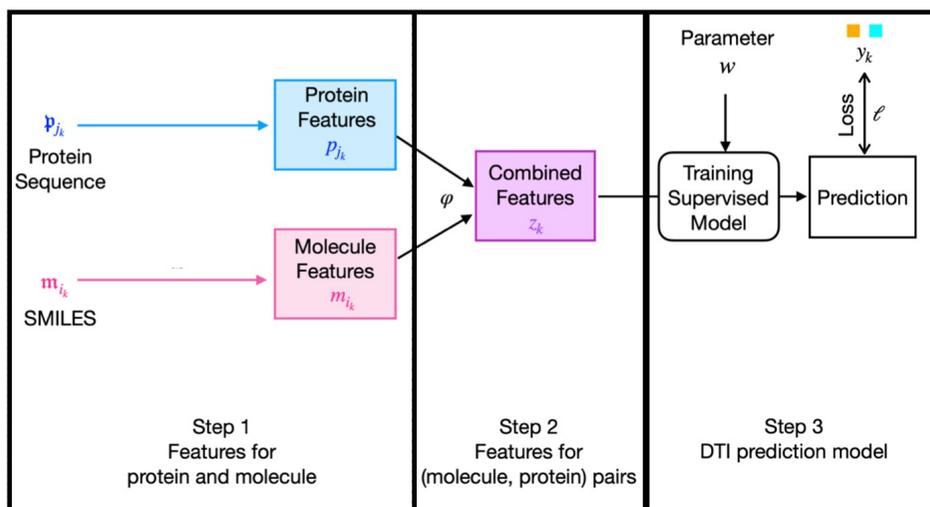
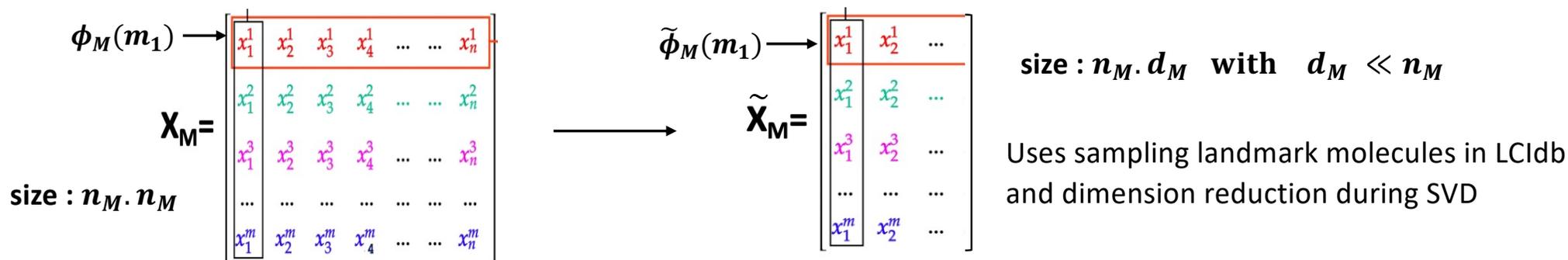
$$\phi_P(p_1) \rightarrow \begin{matrix} \boxed{x_1^1} & x_2^1 & x_3^1 & x_4^1 & \dots & \dots & x_n^1 \\ x_1^2 & x_2^2 & x_3^2 & x_4^2 & \dots & \dots & x_n^2 \\ x_1^3 & x_2^3 & x_3^3 & x_4^3 & \dots & \dots & x_n^3 \\ \dots & \dots & \dots & \dots & \dots & \dots & \dots \\ x_1^m & x_2^m & x_3^m & x_4^m & \dots & \dots & x_n^m \end{matrix}$$

The matrix  $X_P$  is represented as a grid of elements  $x_i^j$ . The first row is highlighted with a red box, and the first column is highlighted with a blue box. The elements are arranged in a grid with rows and columns labeled with indices  $i$  and  $j$ .

Matrix representation  
of the protein data

Compute molecule mapping function  $\phi_M$  : not feasible because  $X_M$  and  $K_M$  too large (size  $n_M \cdot n_M = (271.180)^2$  in LCIdb)

Komet efficiently computes  $\tilde{X}_M$  of much smaller size, such that :  $K_M \approx \tilde{X}_M \cdot \tilde{X}_M^T$



**Step 1:** compute  $\phi_P(p)$  and  $\tilde{\phi}_M(m)$

**Steps 2 and 3:** “classical” SVM optimisation computation of hyperplane in space  $\phi((m, p)) = \tilde{\phi}_M(m) \otimes \phi_P(p)$

**KOMET efficiently solves the optimisation**

**molecules:** SMILE format, **ECFP4 fingerprints** = 1024-bit binary vector,  $K_M = \text{Tanimoto kernel}$   
 $m_M = 3000$  and  $d_M = 1000$

**proteins:** aa sequence,  $K_P$  **LA-kernel** (local alignment).

## KOMET prediction performances and computational efficiency (LCIdb)

Table 4: Comparison of AUPR scores on large-sized datasets, in 5-fold cross-validation.

Dataset	Komet	ConPLex	MolTrans	RF with concatenated features
LCIdb	<b>0.9925±0.0004</b>	0.9783±0.0008	0.9721±0.0011	0.9865±0.0002
Unseen_drugs	<b>0.9944±0.0003</b>	0.9831±0.0009	0.9710±0.0004	0.9829±0.0006
Unseen_targets	<b>0.8952±0.0186</b>	0.8780±0.0223	0.5987±0.0131	0.6886±0.0232
Orphan	<b>0.8671±0.0075</b>	0.8175±.0130	0.5455±0.0004	0.5961±0.0070

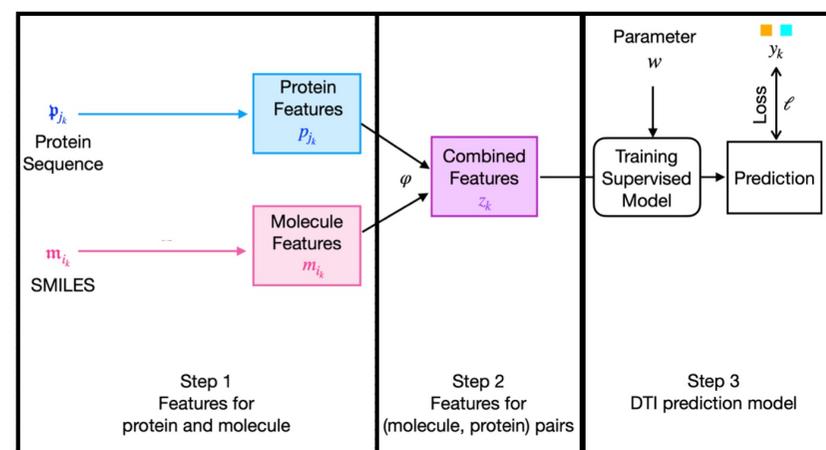
Table 5: Comparison of training time for the considered algorithms.

	Komet	ConPLex	MolTrans	RF with concatenated features
LCIdb	<b>15s</b>	907.3s	69838s	4391s
Unseen_drugs	<b>15s</b>	1734s	68400s	4213s
Unseen_targets	<b>15s</b>	888s	64800s	4100s
Orphan	<b>8s</b>	1329s	25200s	1297s

2 CPUs and 1 NVIDIA A40 GPU  
with 48 GB of memory

Deep-learning algorithms

*Guichaoua et al, JCIIM 2024*



combine with NN

**KOMET also outperforms deep-learning approaches on medium size datasets (BIOSNAP, BindingDB, DrugBank)**

## Target prediction for CFTR modulators with ML chemogenomics

For **Elexacaftor**, Tezacaftor, and Ivacaftor (ETI= Trikafta™)

No high confidence predictions for Ivacaftor and Tezacaftor

For **Elexacaftor**: various **kinases** with high DTI probability scores (> 0.8)

50 kinases experimentally tested.

Modest but significant anti-kinase activity ( $\mu\text{M}$  range) for: **SYK**, **GSK3B** (*in vitro* tests)

Confirmed by *in vivo* functional tests

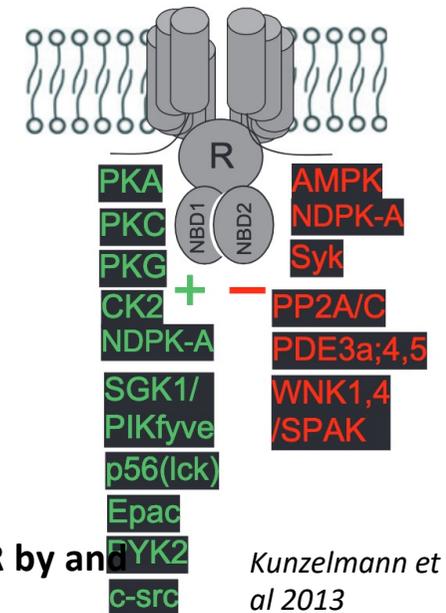
### **SYK and GSK3B appear highly relevant in the context of CF**

- CFTR phosphorylation by **SYK** reduces its stability at the membrane (recycling)

- **GSK3B** inhibitor Kenpaullone improves maturation and processing of F508del-CFTR

(Trzcinska 2012)

**Moderate inhibition of SYK and GSK3B by Elexacaftor may contribute to rescue of mutated CFTR by and increase Chloride conductance**

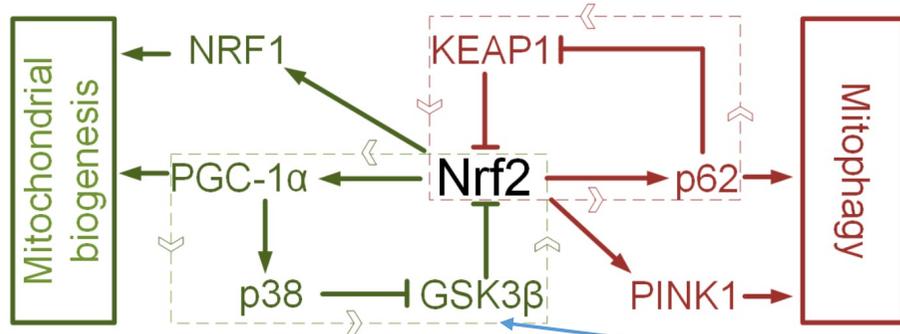


## Further relevance of GSK3B new therapeutic target in CF

### ETI combination improves mitochondrial fitness in CF airway epithelial cells:

- Increases mitochondrial mass, restores Oxidative Phosphorylation activity
- Increases level of PGC1- $\alpha$  (major mitochondrial biogenesis regulator)
- Increases NRF2 levels (TF involved in response to oxidative stress)

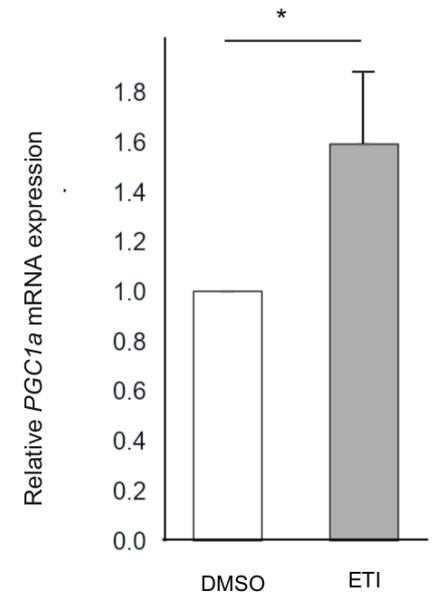
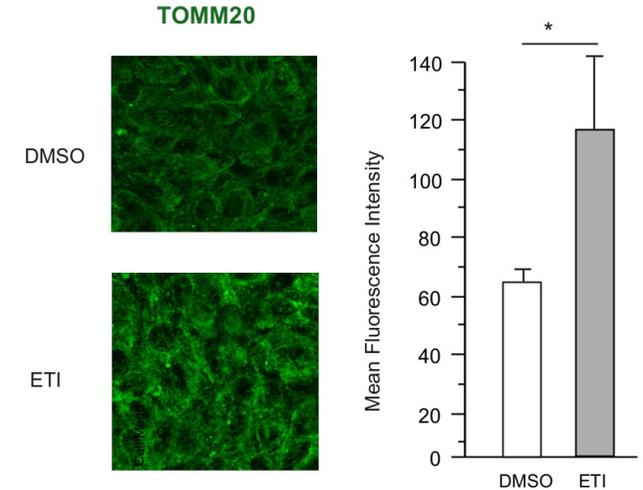
From Guerev *et al*, Brain sciences 2020



**Inhibition of GSK3B by Elexacaftor may explain metabolic benefit of ETI**

Specific inhibitors of GSK3B could improve CF mitochondrial fitness

Unpublished results



## Further relevance of SYK as therapeutic target in CF

Systems Biology approach of CF (Najm et al 2024):

**Rationale: Absence of CFTR propagates further molecular deregulation via its interaction partners leading to various cellular phenotypes that explain diversity of CF symptoms**

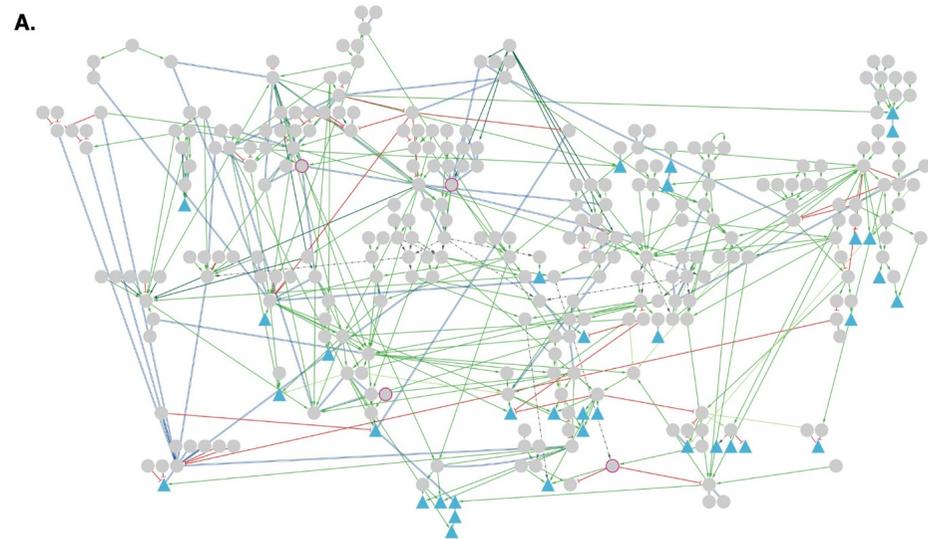
Goal: Build a signalling network that recapitulates CF signalling deregulations

- better understanding of the disease
- identify new therapeutic target whose modulation may modulate CF phenotypes

Network built based on transcriptomic data of CF airway  
Epithelial cells vs WT counterparts

317 nodes, 529 interactions

*Najm et al, BMC Genomics 2024*



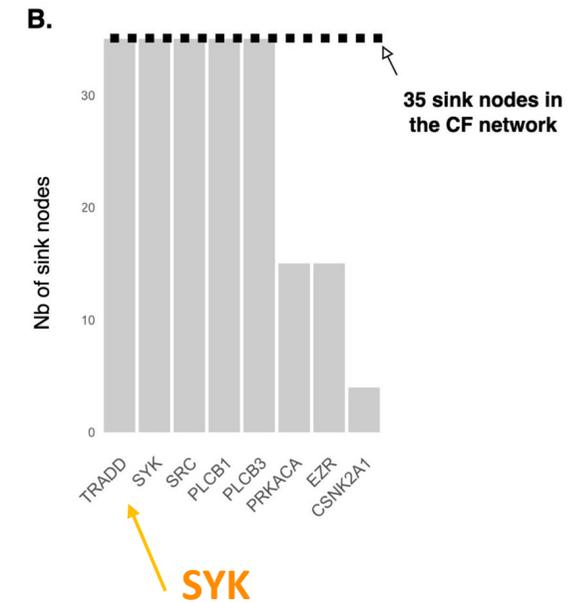
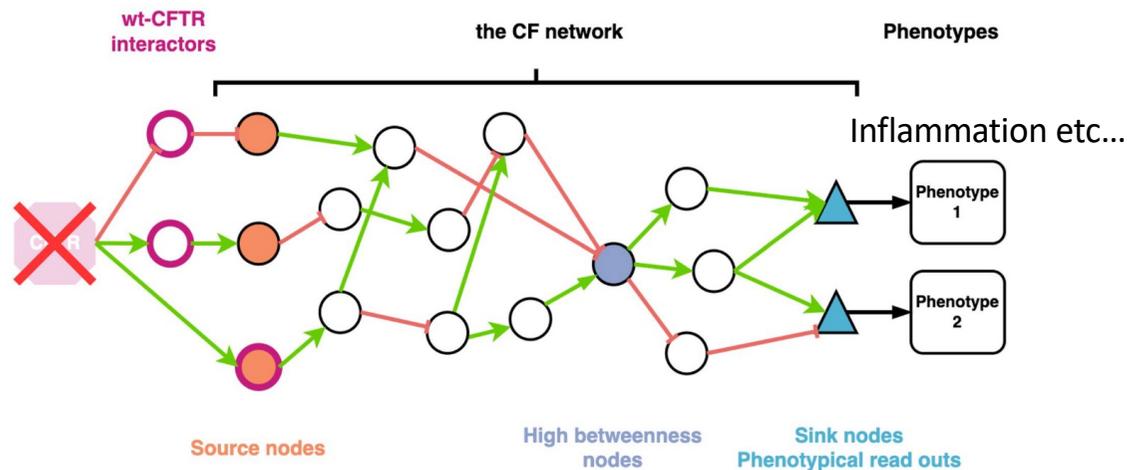
# Identify therapeutic targets based on topological analysis of the network

”remarkable” nodes

**Source nodes** (initiate propagation of deregulations)

**Hubs** (receive and propagate deregulations)

**Output nodes** (modulate CF phenotypes)



Among source nodes, SYK is connected to the 35 sink nodes.

Moderate inhibition of SYK by Elexacaftor may contribute to reduction of inflammation by ETI

**Specific inhibition of SYK may modulate all CF phenotypes in the network**

*Najm et al, BMC Genomics 2024*

## Ongoing work

Further experimental validations of SYK, GSK3B on CF airway respiratory cells using readouts related to CF phenotypes (other than chloride conductance)

**Effect of available specific inhibitors of SYK and GSK3B on :**

- inflammation markers
- Trans Epithelial Electrical Resistance (TEER)
- wound repair (scratch tests)

**Longer term:**

SYK : an interesting target for numerous inflammatory diseases. Fostamatinib FDA approved, others in development

No specific inhibitor of GSK3B approved yet, but several are in development

→ **"mutation agnostic" therapies** to modulate deleterious CF symptoms

## Acknowledgements

**CBIO:** Mathieu Najm, Gwenn Guichaoua, Philippe Pinel, Victor Laigle, Chloé-Agathe Azencott



**INEM, U1151:** Isabelle Sermet-Gaudelus, Mairead Aubert



**Institut Curie, U1331:** Loradana Martignetti, Laurence Calzone



**Financial Support:** La Fondation Dassault Systèmes, Vaincre la Mucoviscidose, Fondation pour la Recherche Médicale, La Fondation Maladies Rares

Datasets	Molecules	Proteins	Positive DTIs	Negative DTIs
BIOSNAP	4,510	2,181	13,836	(13,647 random)
Unseen_drugs			13,836	(13,647 random)
Unseen_targets			13,836	(13,647 random)
BindingDB	7,161	1,254	9,166	23,435
DrugBank	4,813	2,507	13,715	(13,715 balanced)
DrugBank (Ext)	4,257	1,216	10,838	(10,838 balanced)
LCIdb	271,180	2,060	396,798	7,965 (+ 388,833 balanced)
Unseen_drugs	271,180	2,060	396,798	7,965 (+ 388,833 balanced)
Unseen_targets	271,180	2,060	396,798	7,965 (+ 388,833 balanced)
Orphan	191,901	2,060	208,041	7,965 (+ 200,076 balanced)

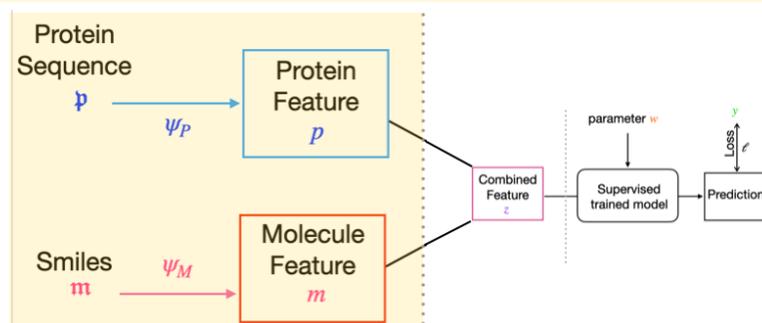
**LCIdb**

**30 times larger than BIOSNAP  
and DrugBank**

**43 times larger than BindingDB**

*Guichaoua et al, JCIM 2024*

# Impact of molecule and protein features choice



## Comparison of fixed and pre-trained learned features

Performance in the more difficult scenario: LCldb\_Orphan Dataset

AUPR		Protein feature	
		Feature map of LA kernel	ESM2 [Lin et al, 2022]
Molecule Feature	Feature map of Tanimoto kernel on ECFP4	0.897 🚀	0.864
	ECFP4	0.893	0.866
	Dgl-lifesci [Li et al, 2021] (GNN supervised contextpred)	0.887	0.858

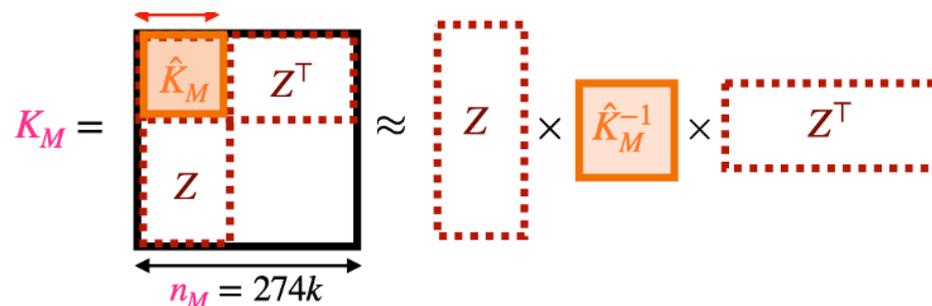
## Fixed features better than DL features for this specific problem

Drug-like molecules and human druggable proteins

**Molecule kernel**  $k_M(\mathbf{m}, \mathbf{m}') = \langle \psi_M(\mathbf{m}), \psi_M(\mathbf{m}') \rangle$   $K_M = X_M X_M^\top$

**Nyström approximation**

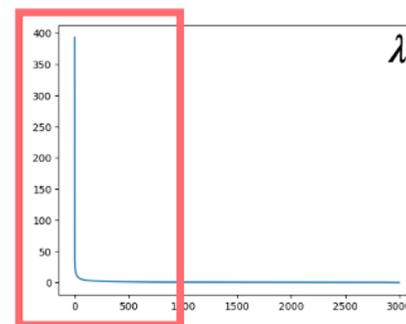
$m_M$  random landmarks molecules



**Singular value decomposition (SVD)**

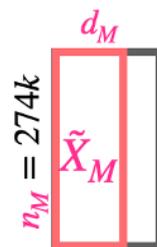
$\hat{K}_M = U \text{diag}(\lambda) U^\top$

$K_M \approx X_M X_M^\top$  where  $X_M = Z U \text{diag}(1/\sqrt{\lambda})$

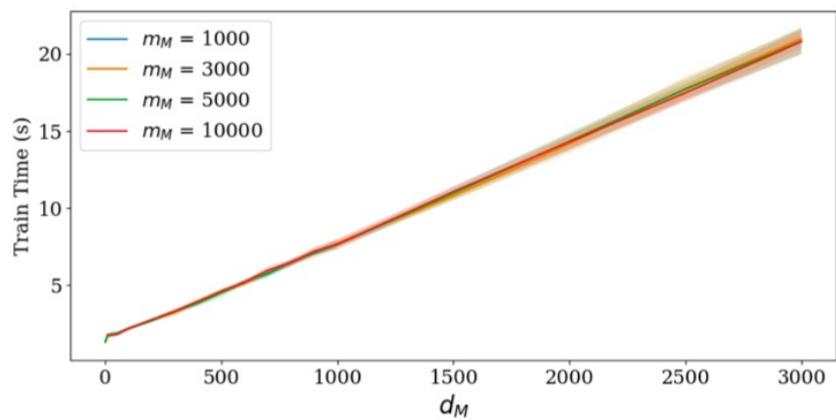
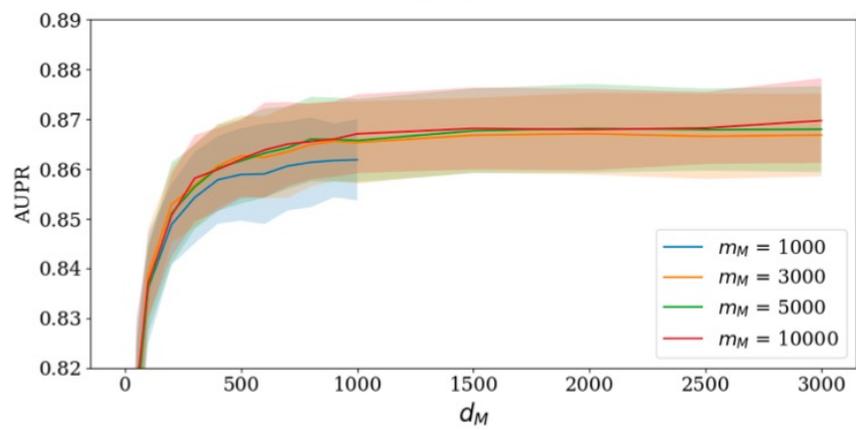


**Reduction dimension**

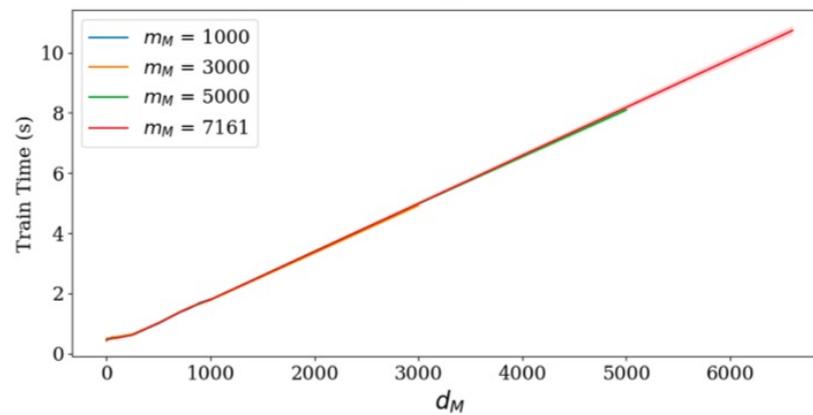
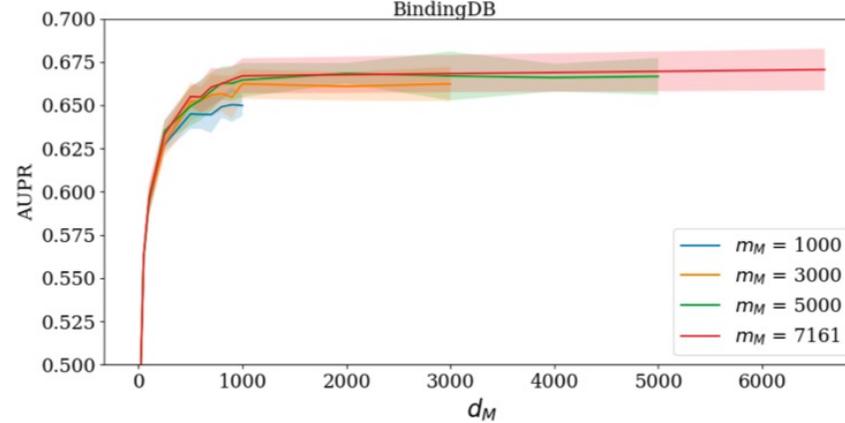
$X_M \approx \tilde{X}_M$  where  $\tilde{X}_M \in \mathbb{R}^{n_M \times d_M}$



LCIdb\_Orphan



BindingDB



## KOMET Prediction performances:

Table 4: Comparison of AUPR scores on large-sized datasets, in 5-fold cross-validation.

Dataset	Komet	ConPLex	MolTrans	RF with concatenated features
LCIdb	<b>0.9925±0.0004</b>	0.9783±0.0008	0.9721±0.0011	0.9865±0.0002
Unseen_drugs	<b>0.9944±0.0003</b>	0.9831±0.0009	0.9710±0.0004	0.9829±0.0006
Unseen_targets	<b>0.8952±0.0186</b>	0.8780±0.0223	0.5987±0.0131	0.6886±0.0232
Orphan	<b>0.8671±0.0075</b>	0.8175±0.0130	0.5455±0.0004	0.5961±0.0070

On large size datasets (LCIdb)

Table 5: Comparison of training time for the considered algorithms.

	Komet	ConPLex	MolTrans	RF with concatenated features
LCIdb	<b>15s</b>	907.3s	69838s	4391s
Unseen_drugs	<b>15s</b>	1734s	68400s	4213s
Unseen_targets	<b>15s</b>	888s	64800s	4100s
Orphan	<b>8s</b>	1329s	25200s	1297s

On medium size datasets

Dataset	Komet	ConPLex	MolTrans	RF with concatenated features
BIOSNAP	<b>0.9429±0.0008</b>	0.9246±0.0037	0.8989±0.0048	0.9121±0.0032
Unseen_drugs	<b>0.8979±0.0051</b>	0.8763±0.0072	0.8547±0.0045	0.8620±0.0100
Unseen_targets	<b>0.8754±0.0099</b>	0.8641±0.0100	0.7058±0.0273	0.8127±0.0116
BindingDB	0.6598±0.0074	<b>0.6765±0.0178</b>	0.6196±0.0150	0.6454±0.0075
DrugBank	<b>0.9400±0.0030</b>	0.8961±0.0070	0.8068±0.0100	0.8018±0.0086

Deep-learning algorithms