# Drug Discovery as a Recommendation Problem

Anna Gogleva, Erik Jansson,
Greet De Baets, Eliseo Papa
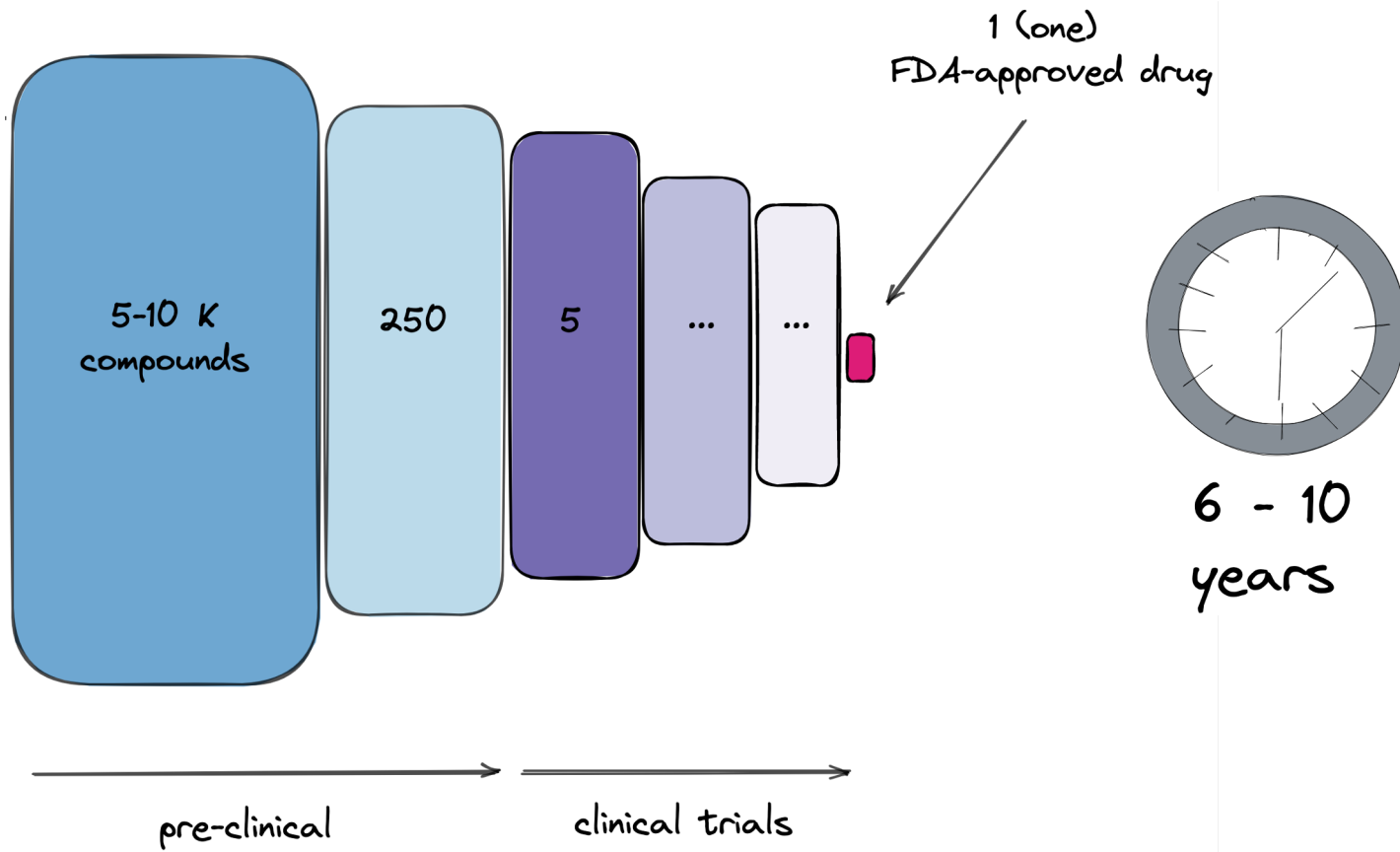
27th September 2021

ACM RecSys'21 Amsterdam
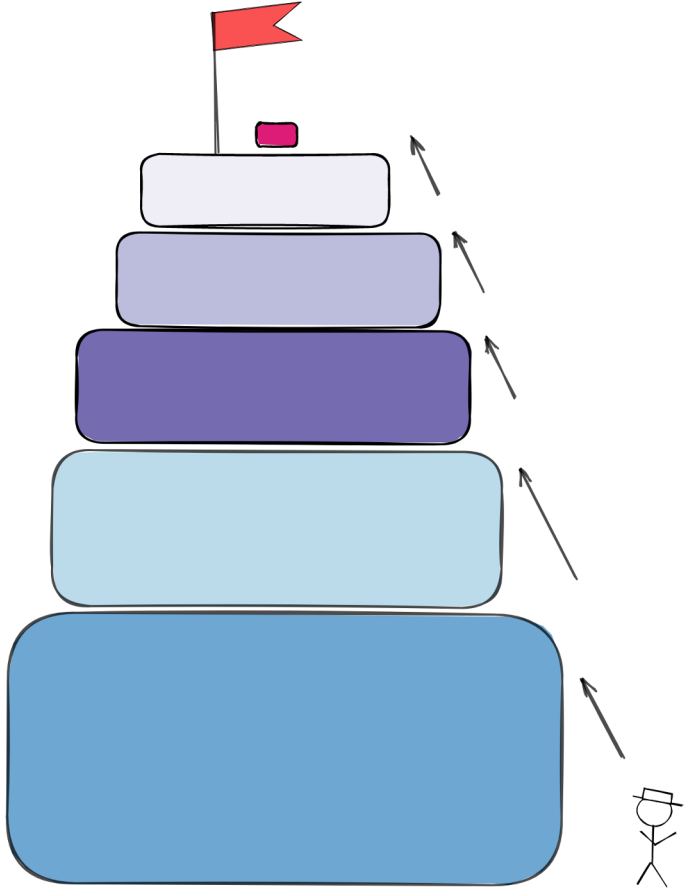
**https://astrazeneca.github.io/recsys21gogleva/**

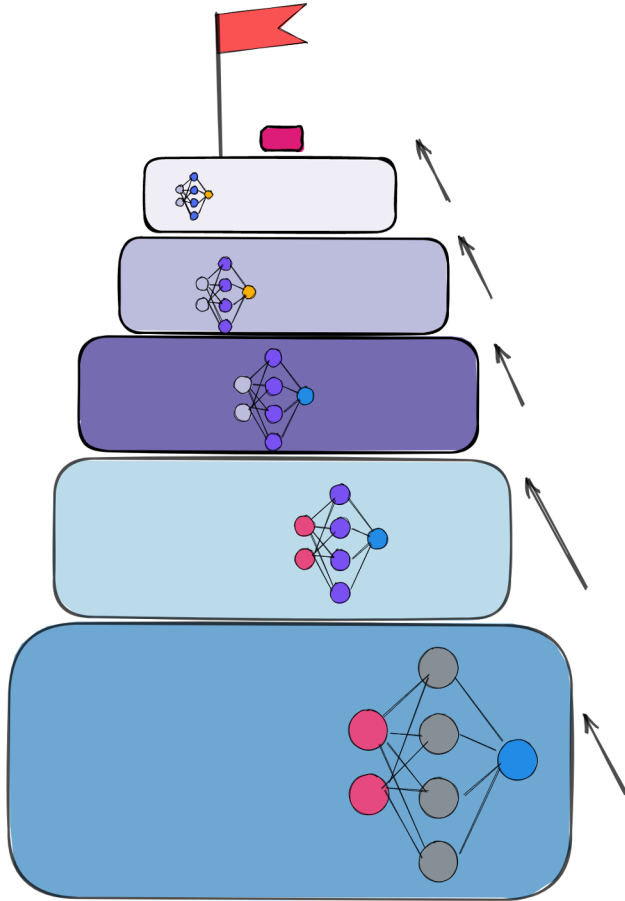# One needs to fail a lot to discover a working drug
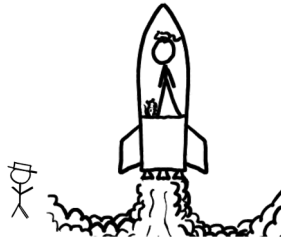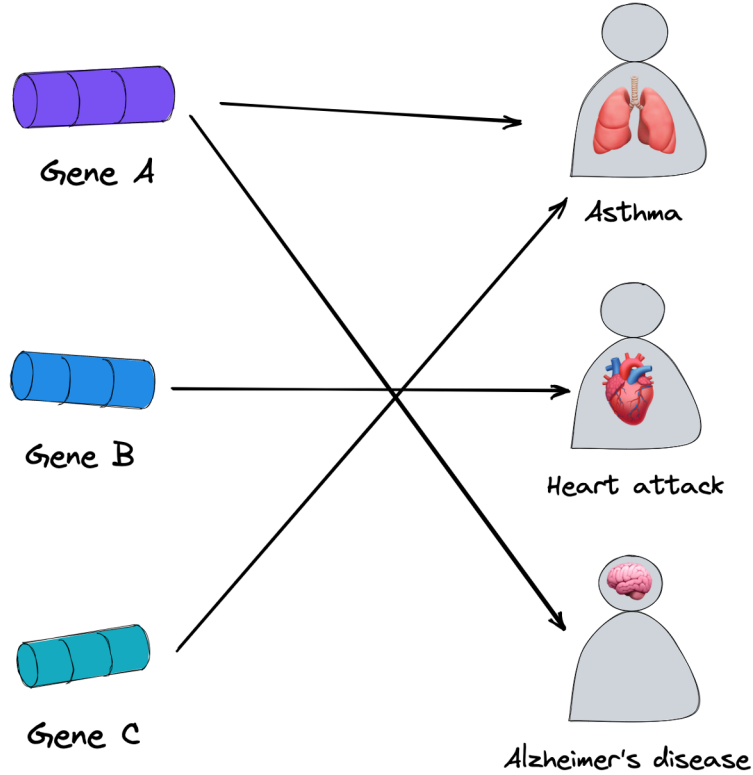
# It is a tall mountain to climb



- How to develop new efficient treatments faster?

- How to make better decisions in the process?

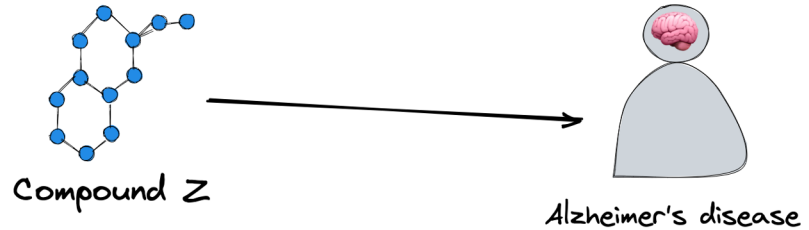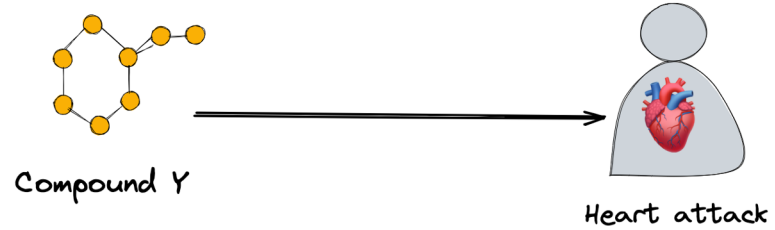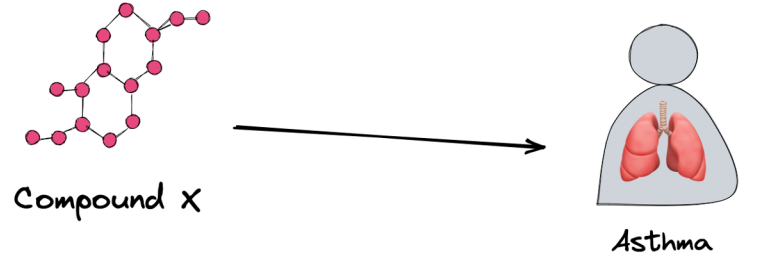# It is a tall mountain to climb



- How to develop new efficient treatments faster?

- How to make better decisions in the process?

- Recommendation systems can help in multiple places

# Recommendation problems in drug discovery
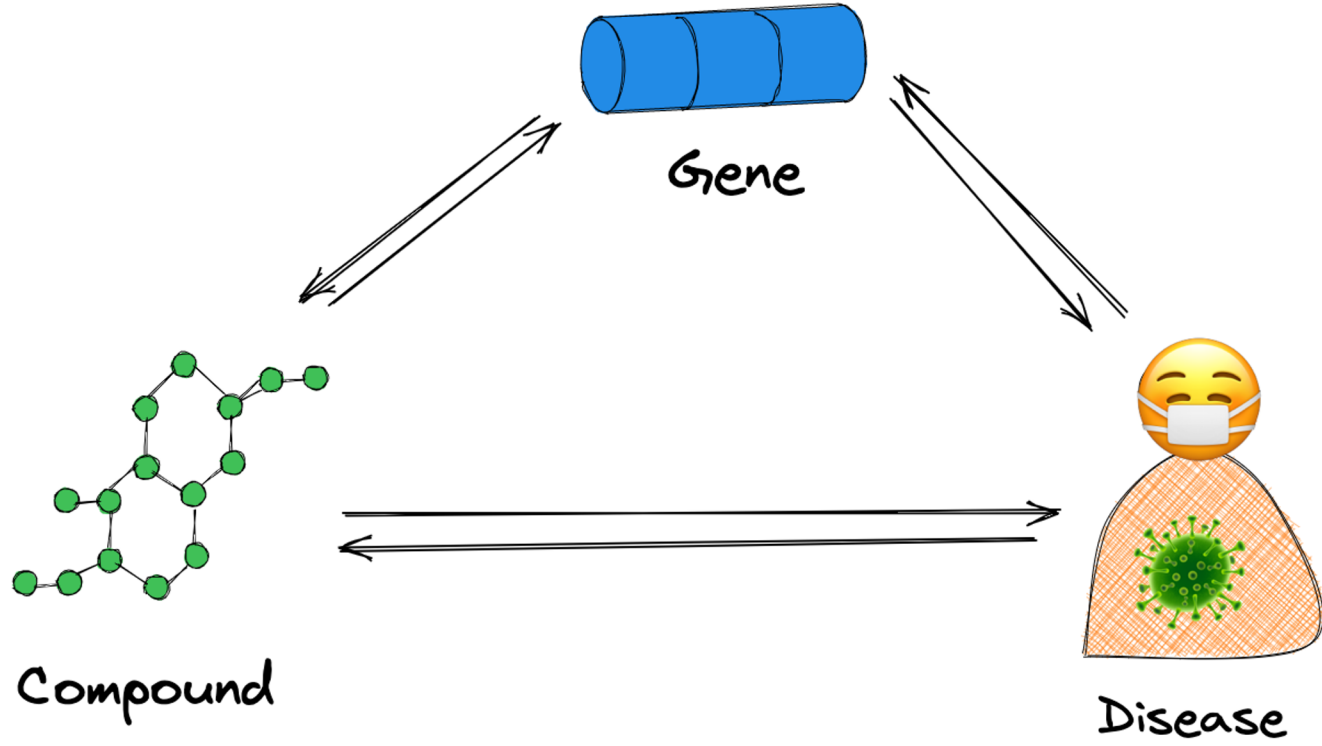


find a gene causing a disease

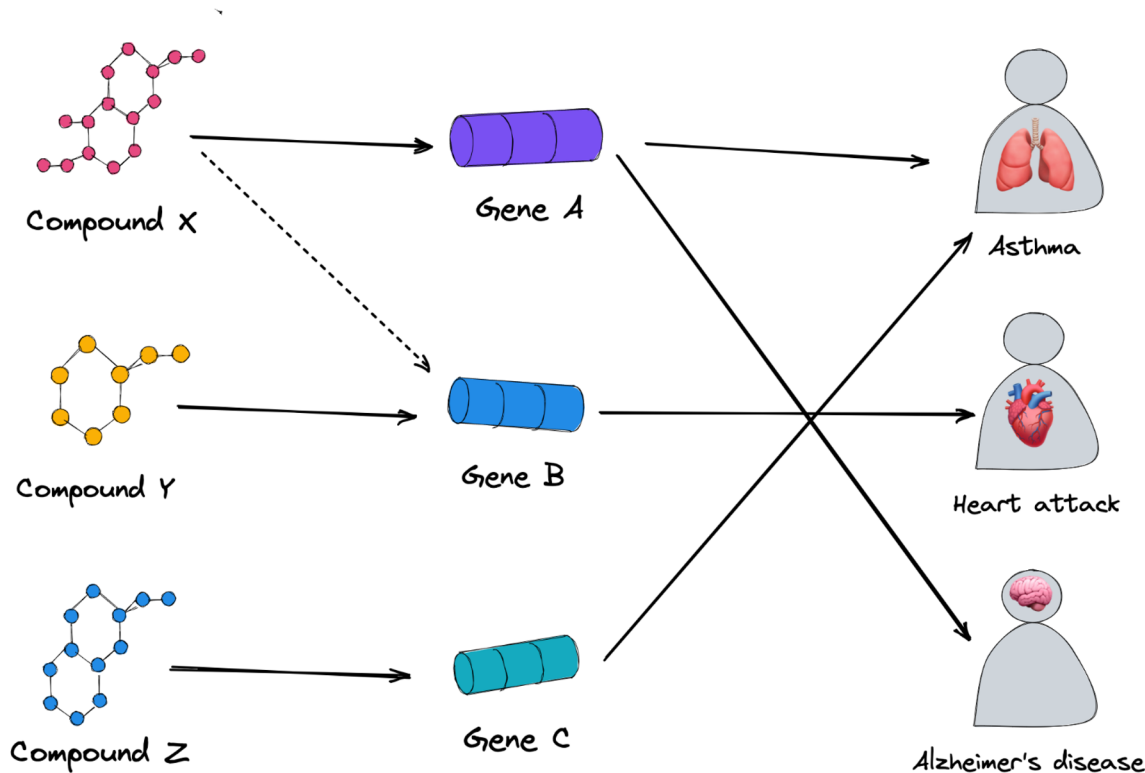match a drug with a disease

# Drugs, genes, diseases

# It gets complex very fast

Millions of compounds
Billions possible theoretically

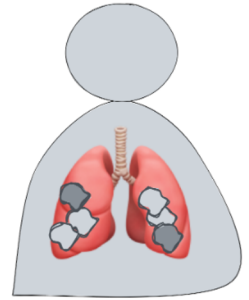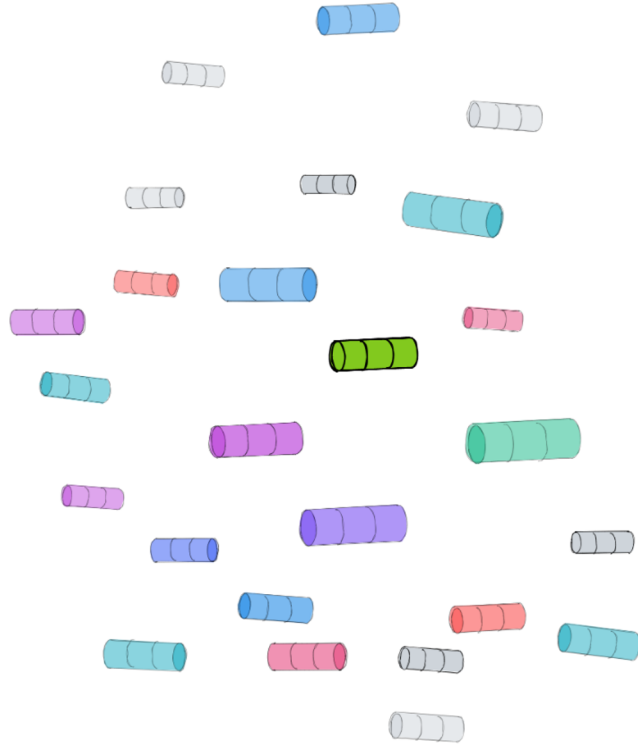25-30 K genes,
80 K functional elements

~10 K diseases

# It is rarely just a single gene

- 25-30K human genes

- everything interacts with everything, each gene is a suspect
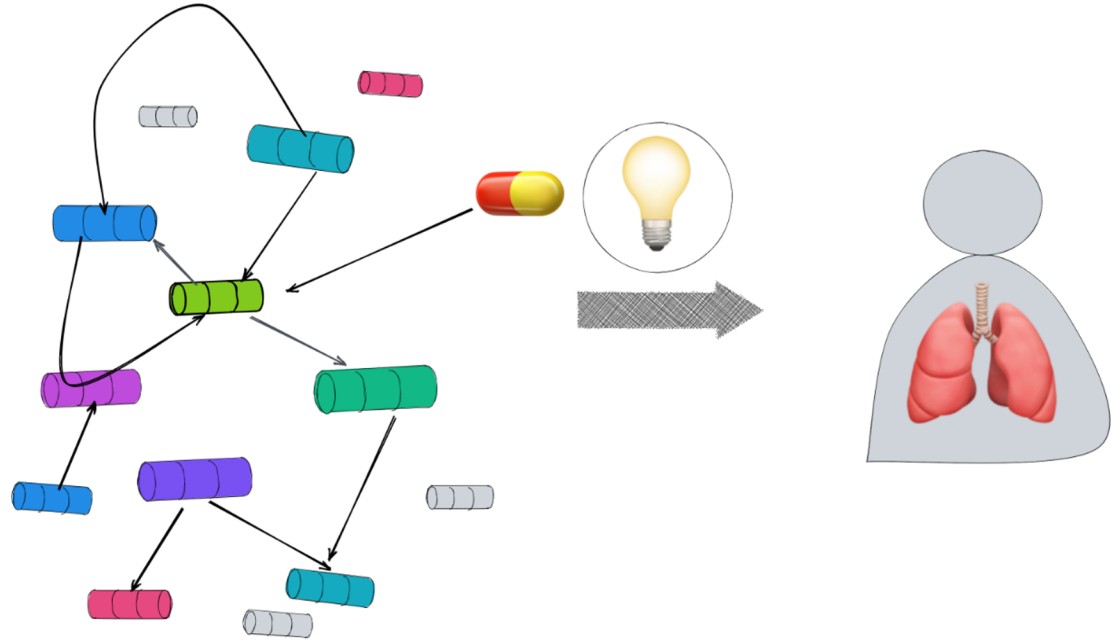


a disease

# Find a molecular network behind a disease



1. disease ~ a molecular process gone awry

2. find the key molecular process

3. re-route it safely

# Biomedical knowledge is spread across multiple resources

# Graph makes things simpler



- Biomedical information often comes in forms of **networks** and **hierarchies**

- Graph is a convenient way to organise it

- BIKG (our internal knowledge graph): **60+** data sources including - omics and data extracted from the literature

- **11 M nodes, 1 B edges**

- Use graph as a source of context and features for recommenders

Early success story:

graph-based
recommendations

# Applied recommendation problem #1: contextualize experimental data



- Drug resistance in lung cancer

- Occurs in a sub-population of patients

- Resistance landscape is complex

*X Wang*, H Zhang, *X Chen* - *Cancer Drug Resistance, 2019*

# How to help scientist find key genes faster?



- key gene == target
- find a target gene ➡️ develop new treatments to overcome the resistance

# An ideal target

- - -



- ☑ Expression
- ☑ Pathway/complex enrichment
- ☑ Effect size
- ☑ Druggability
- ☑ Mode of action
- ☑ Translation in models
- ☑ Internal assets
- ☑ Bench validation
- ☑ Consistency in assays
- ☑ Clinical relevance
- ☑ Literature support
- ☑ Novelty

...

# An ideal target does not exist

- - -



☑ Expression

☑ Pathway/complex enrichment

☑ Effect size

☑ Druggability

☑ Mode of action

☑ Translation in models

☑ Internal assets

☑ Bench validation

☑ Consistency in assays

☑ Clinical relevance

☑ Literature support

☑ Novelty

...

# Target selection as an optimization problem

# Hybrid feature set: source features from the graph



Betweenness

Node degree

PageRank

Clustering coefficient

Graph embeddings

# Hybrid feature set: combine with clinical features



Betweenness

Node degree

Clinical features

PageRank

Literature support

Druggability

Clustering coefficient

Graph embeddings

Pre-clinical experimental assays

# Approaches



① Compute exact Pareto front

② Evolutionary algorithms

evaluation     selection

mutation     crossover

③ Matrix factorization

Implicit feedback

level n

...

level 1

f1

f2

W     U     V

# SkywalkR, interactive interface

- - -

- select a subset of objectives

- set optimization directions

- explore trade-offs



github.com/AstraZeneca/skywalkR

# Imperfect validation

# Model domain scientist as a black box classifier



(Gogleva et al, biorXiv, 2021)

# Graph-derived features follow clinical in unbiased setting



(Gogleva et al, biorXiv, 2021)

# Annotation by the experts



WWTR1      WW domain containing transcription regulator 1
ENSG00000018408

#Publications of this hit mentioned within the context of 'resistance' and 'EGFR': 0

#Publications of this hit mentioned within the context of 'resistance' and 'NSCLC': 0

for additional evidence behind the gene recommendation please see skywalkR

☐ Known resistance marker                                    1

☐ Novel, but credible hit                                    2

☐ Novel, not credible hit                                    3

☐ Not novel, not credible hit                                4
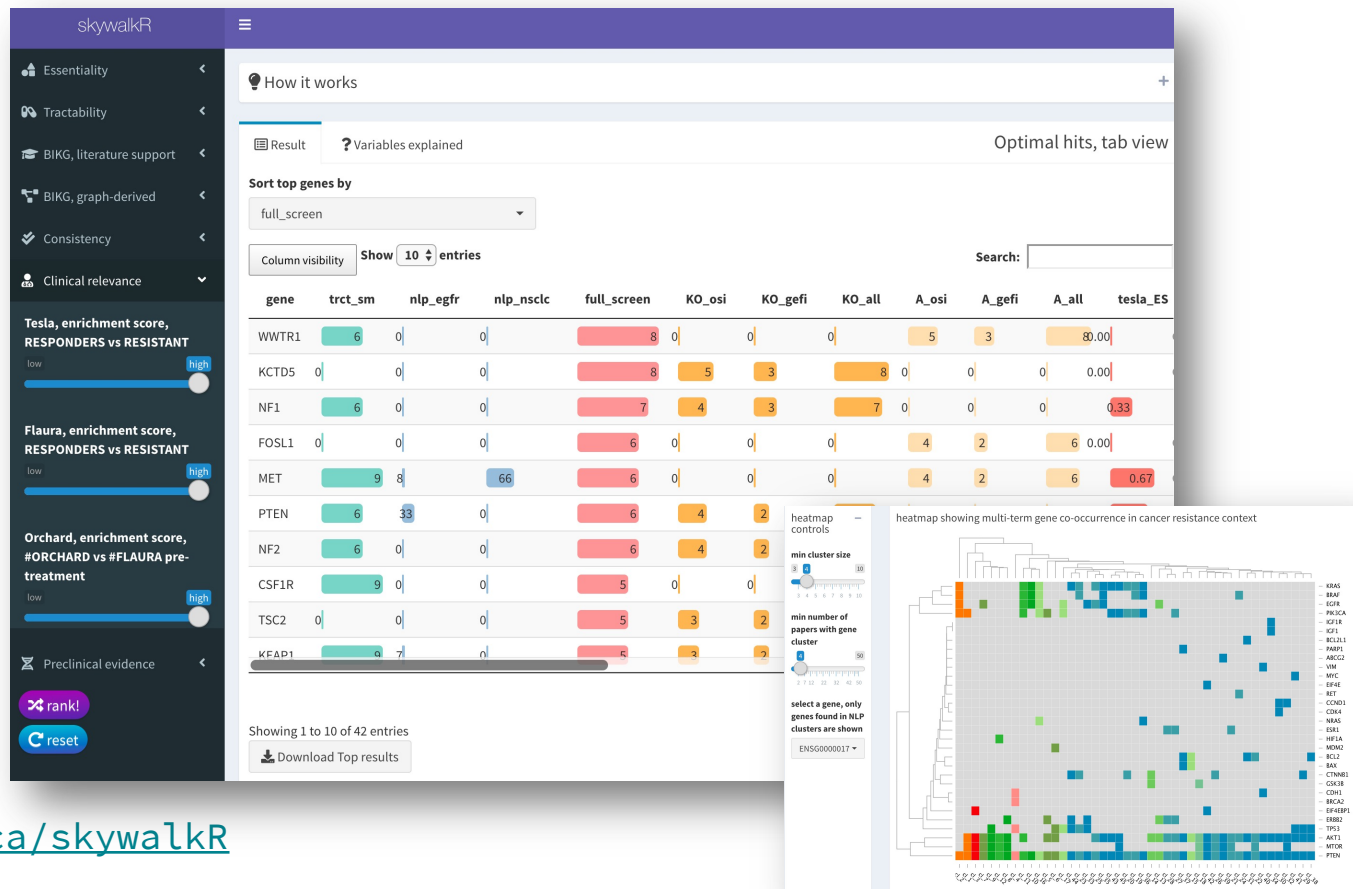
please include any additional details about ongoing experiments for this marker, or if this has been discussed at (pre)TSID.

TASK_NUM: 1   TOTAL_TASKS_NUM: 42

Gene list

prodigy

# Most of recommendations are 'novel & credible'



(Gogleva et al, biorXiv, 2021)

# Experimental validation *in vitro*



- confirmed involvement of 4 recommended genes in drug resistance

- next: test the remaining genes

(Gogleva et al, biorXiv, 2021)

# Imperfect, yet already useful recommendation system



- 🐌 -> 🏍 re-rank lists in seconds, not months
- ⚙️ automated feature generation
- ♻️ approach can be re-used in related problems

# Take home message

___

- Drug discovery is an exciting field for recommender systems

- Relatively simple recommenders can have a lot of impact

- Need for recommenders that can operate in unsupervised or weakly supervised settings

- There are a number of challenges✴️

✴️ Read more in the extended deck:
    **https://astrazeneca.github.io/recsys21gogleva/**

**Collaborate with us!**

anna.gogleva@astrazeneca.com

bikg@astrazeneca.com

BIKG

# Translating recommendation approaches to biomedical field: a few complications

# Biological entities are complex



complex gene interactions

diseases can be poorly defined

gene A     causes     disease X

multiple disease mechanisms

diseases can be heterogeneous

# Validation is slow and expensive



validation requiers experiments

in vitro

in vivo

in clinic

feedback

needs seconds

can take years

# Implicit & explicit feedback is scarce

# Team of experts rather than a single user makes decisions



feedback

needs seconds

experts

in vitro

in vivo

in clinic

can take years

a team of experts makes decisions

each expert has their own bias

# Previous literature biases users decisions



Moderately studied

Dark matter

Papers

genes

Nature

filter bubble:

a small number of well studied genes tends to get the credit

'dark matter' of human genome remains under-explored

# Ground truths are rare and context-specific

gene A    interacts    gene B
           with

same time:
- disease stage
- developmental stage

same place:
- tissue
- organ

same genetic background

😭 there is a lot of data out there,
    but never the data you need to train your model

# Portfolio problem vs single choice:
## continuously optimize based on constantly changing evidence

# Supplementary: supervised recommendations

# Can we learn from previous drug trials?

___

- Thousands of clinical trials preclinical experiments (internal + external)

- Idea: use data on previous (potential) targets as training data for a supervised model



Pre-clinical      Clinical trials

# Can we learn from previous drug trials?

———

- Represent genes with experimentally derived and KG-derived features
  - Experimental - activity in certain bio processes
  - KG-derived - graph distances, embedding distances, etc. etc.

# Can we learn from previous drug trials?

———

- Train a supervised ranking model (LightGBM) with randomly sampled targets as negatives and clinically promising targets as positives

# Human-Model trust

\_ \_ \_

- We need biologist's to sign off on our model's recommendations
- For that, we need their trust
    - NDCG or other "ML" metrics mean nothing to a biologist
    - Biologists expect certain genes as a sanity-check

# Human-Model trust

—  —  —

"I would expect to see Gene X in your
recommendations - otherwise we have a
problem"

# Human-Model trust

---

"Yup the model is recognizing Gene X as a promising gene target!"

# Human-Model trust

___

"How do I know the model isn't just regurgitating what I've told you?"

# Human-Model trust

———
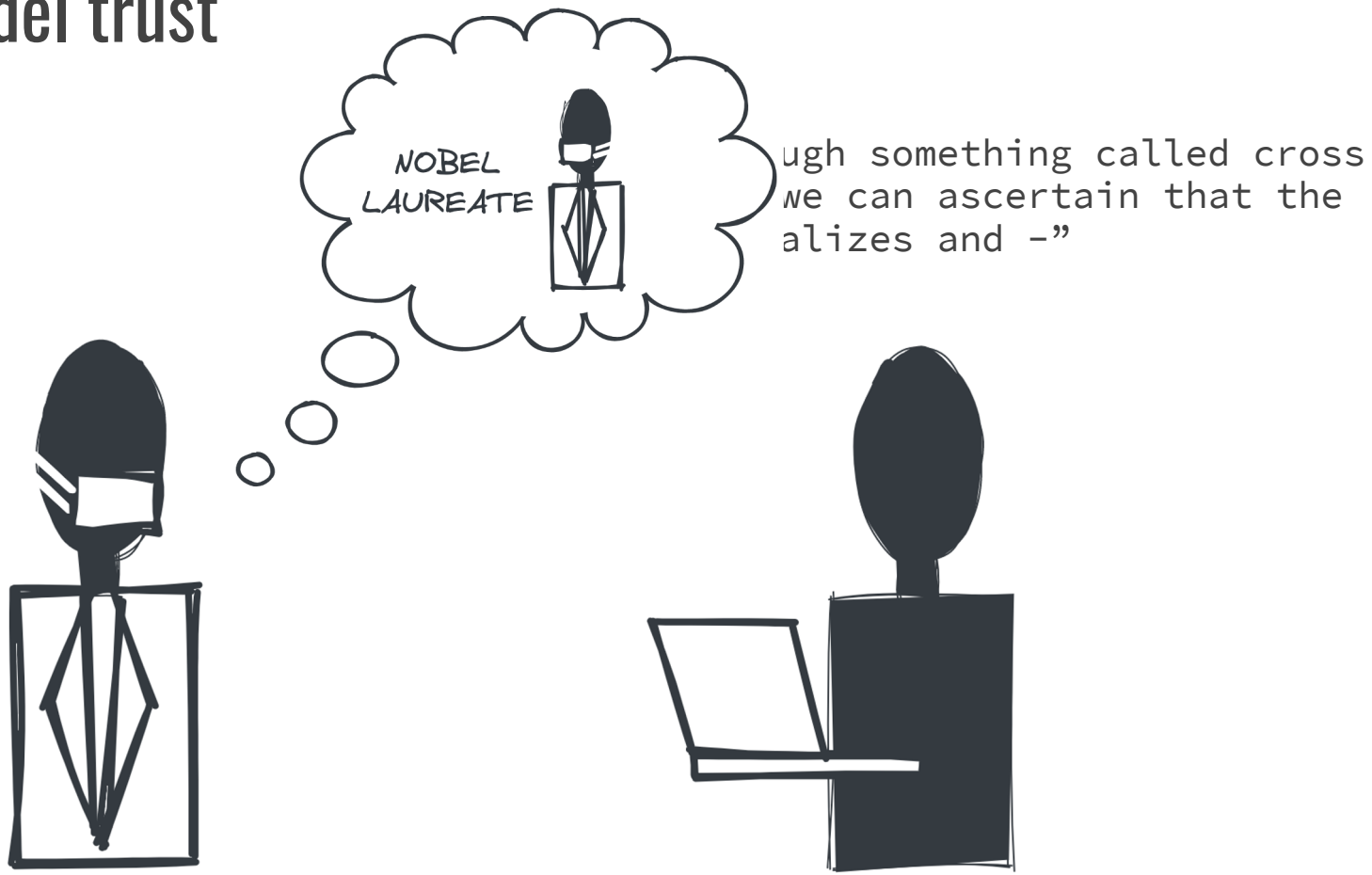
"Well, through something called cross validation we can ascertain that the model generalizes and -"

# Human-Model trust
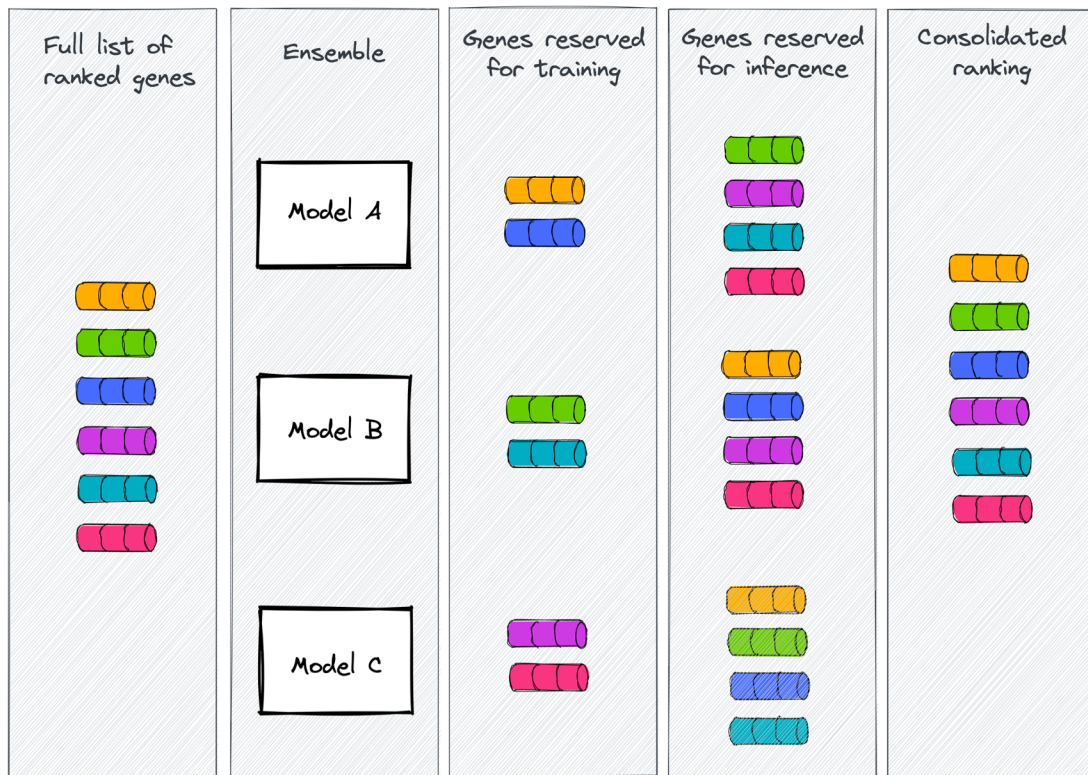
# Human-Model trust

———

- Problem: human genome is finite
  (Since we rank the full genome, the training set will exist somewhere in the
  final model output)

- How can we guarantee that no "regurgitation" is happening
  during inference?

# "Honest" Ensembling

---

- Training data is split among an ensemble of models
- If a gene has been seen by a model during training - this model can't rank its target-aptitude during inference

# "Honest" Ensembling

# Jury is still out

———

- Training data: genes that have previously been found promising in COPD (Chronic obstructive pulmonary disease)
- After ranking:
  - Take the top ~200 genes
  - Filter for known involvement in a number of interesting molecular processes
  - Bring to biologists for manual quality control
- => 29 potential gene targets are now in experimental validation