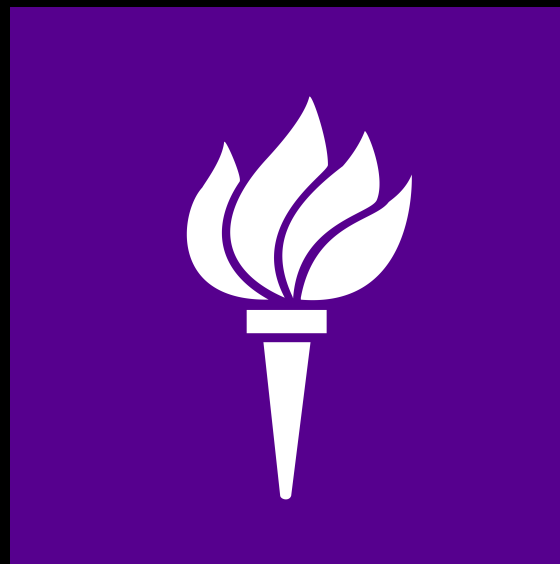
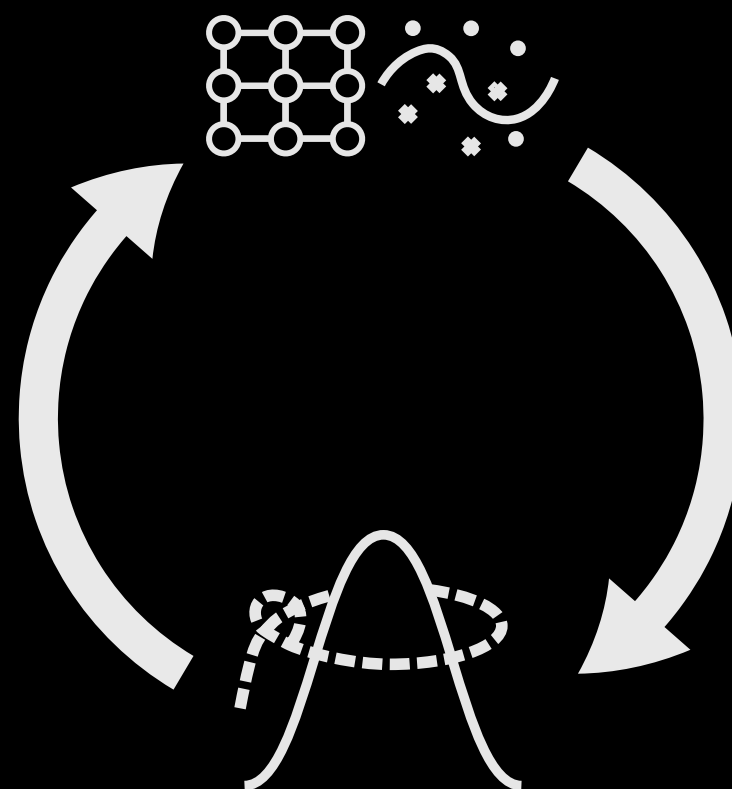
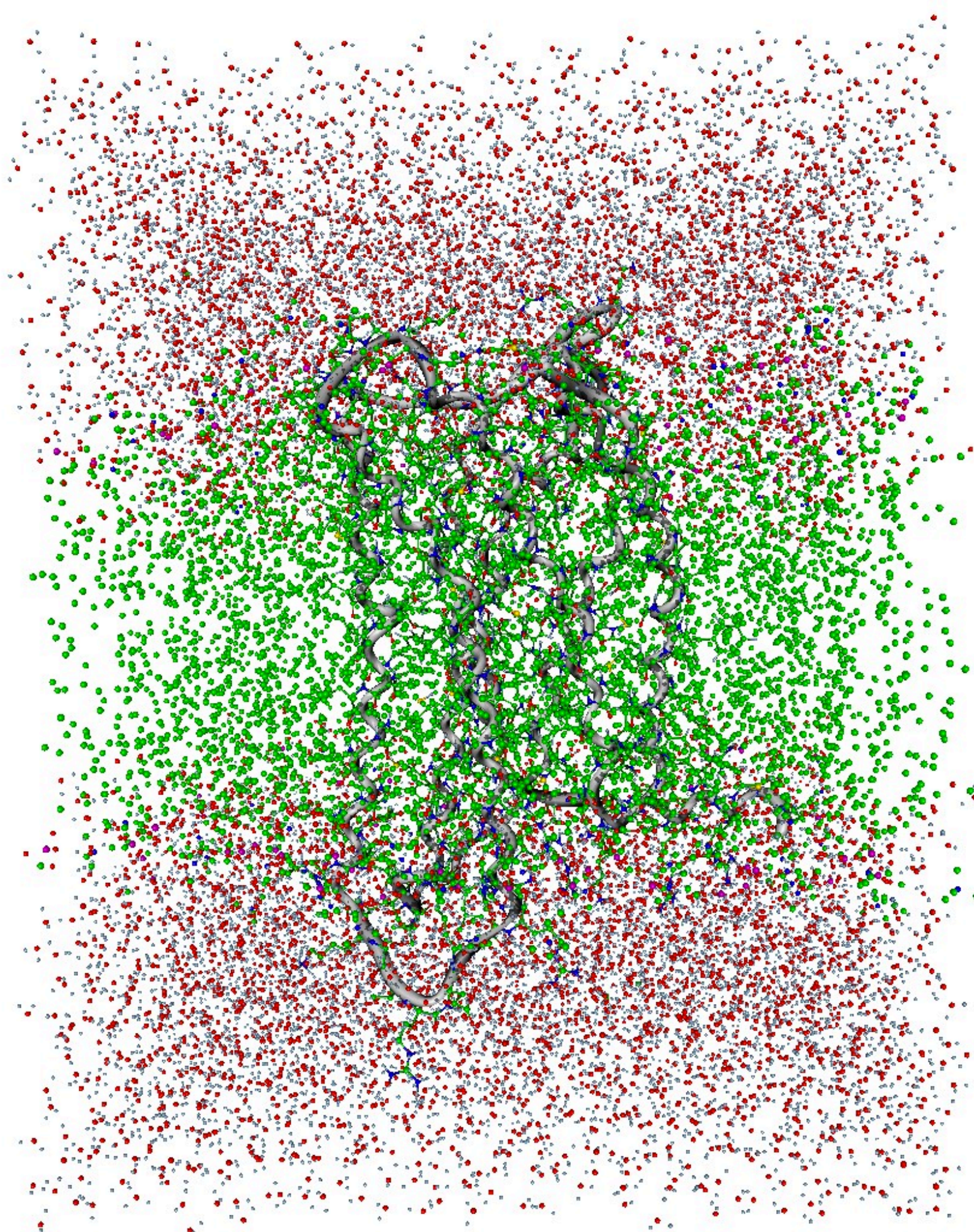
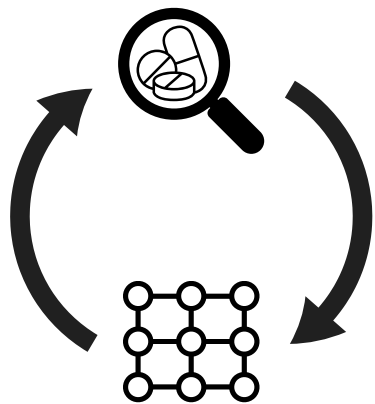


Probabilistic Modeling of Structure in Science: Statistical Physics to Recommender Systems

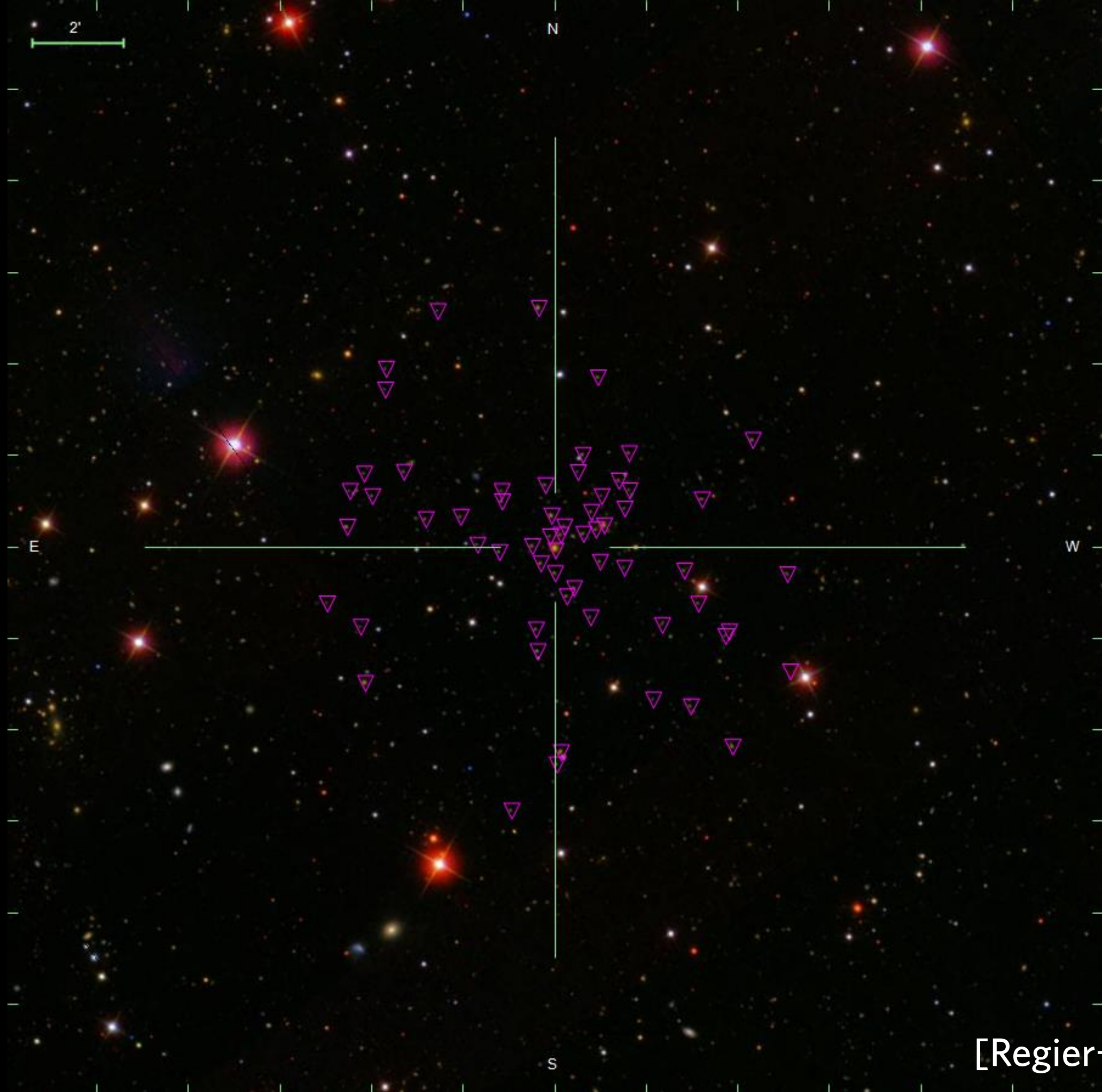
Jaan Altosaar



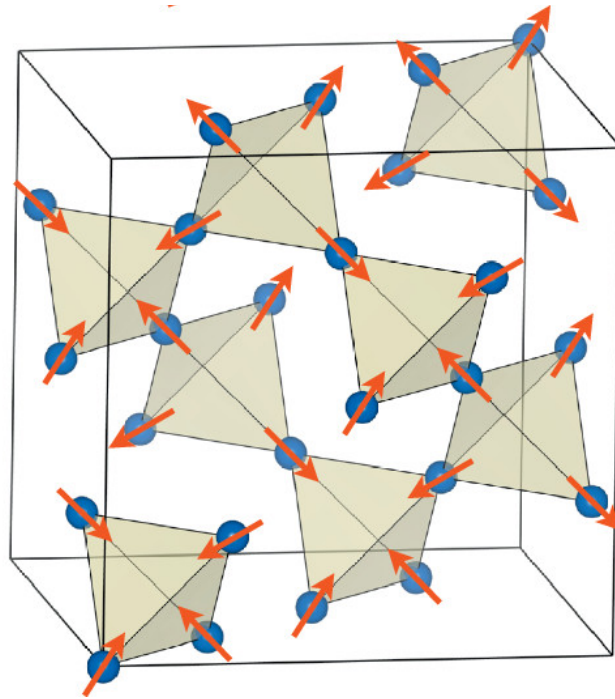




[Klimm+ 2018]

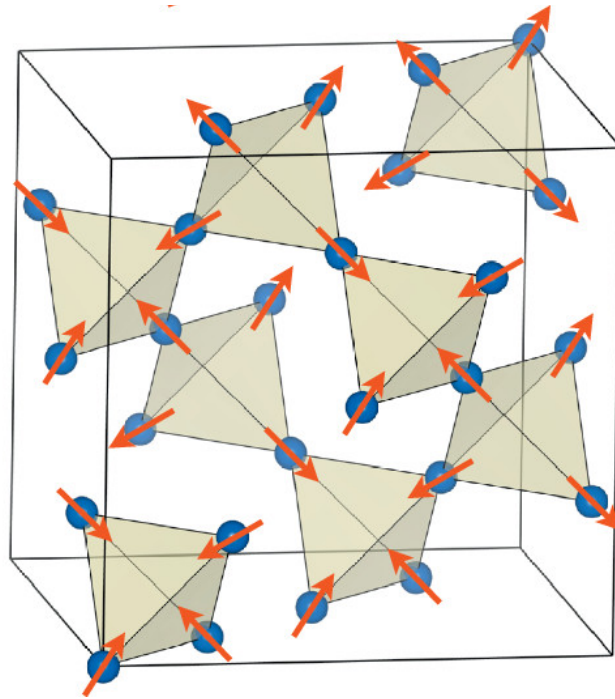


[Regier+ 2018]



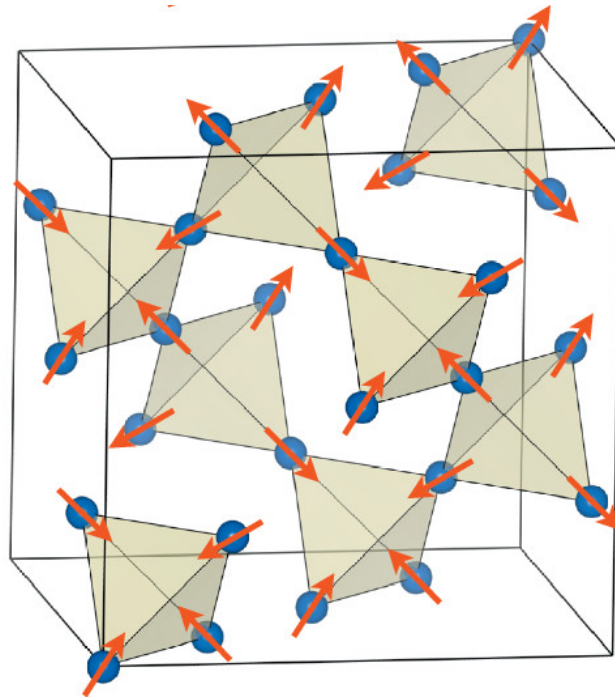
$$C_v = f(\mathcal{Z})$$

$$p(\mathbf{z}) = \frac{e^{-\beta H(\mathbf{z})}}{\mathcal{Z}}$$



$$p(\mathbf{z}) = \frac{e^{-\beta H(\mathbf{z})}}{\mathcal{Z}}$$

$$\mathcal{Z} = \sum_i \sum_{\mathbf{z}_i = \pm 1} e^{-\beta H(\mathbf{z}_i)}$$



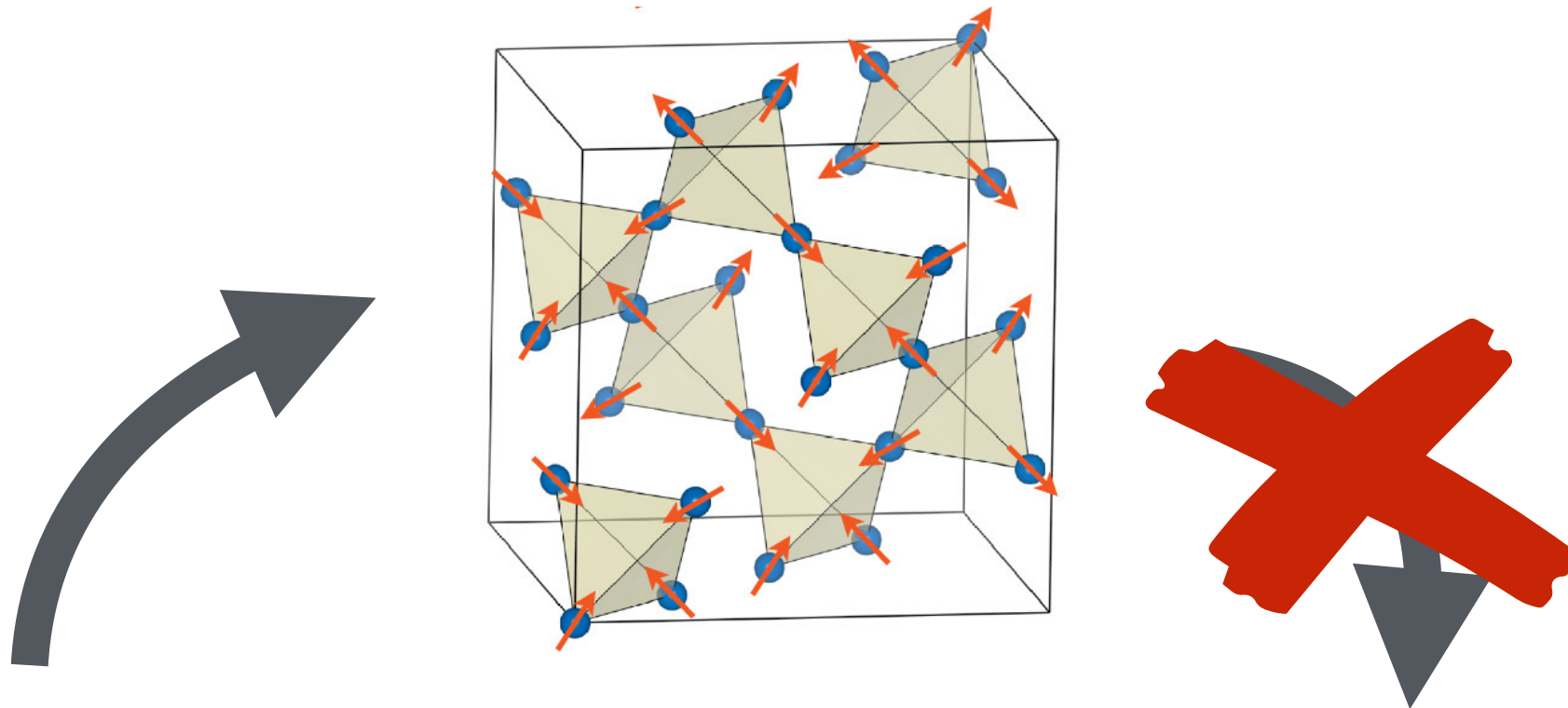
- Markov Chain Monte Carlo

- Approximate $f(\mathcal{Z})$

[Neal 1993; Hastings 1970]

$$p(\mathbf{z}) = \frac{e^{-\beta H(\mathbf{z})}}{\mathcal{Z}}$$





- Variational Inference

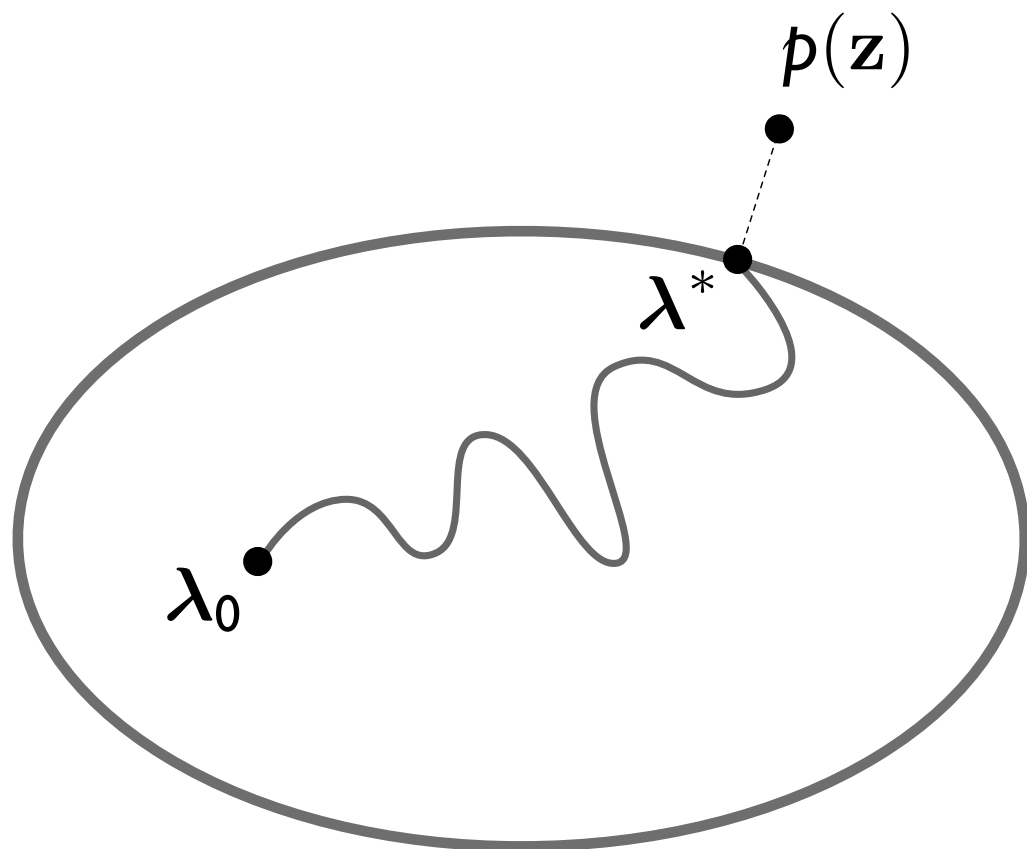
- Learn approximation $q(\mathbf{z})$

[Blei+ 2016; Peterson+ 1987]

$$p(\mathbf{z}) = \frac{e^{-\beta H(\mathbf{z})}}{\mathcal{Z}}$$

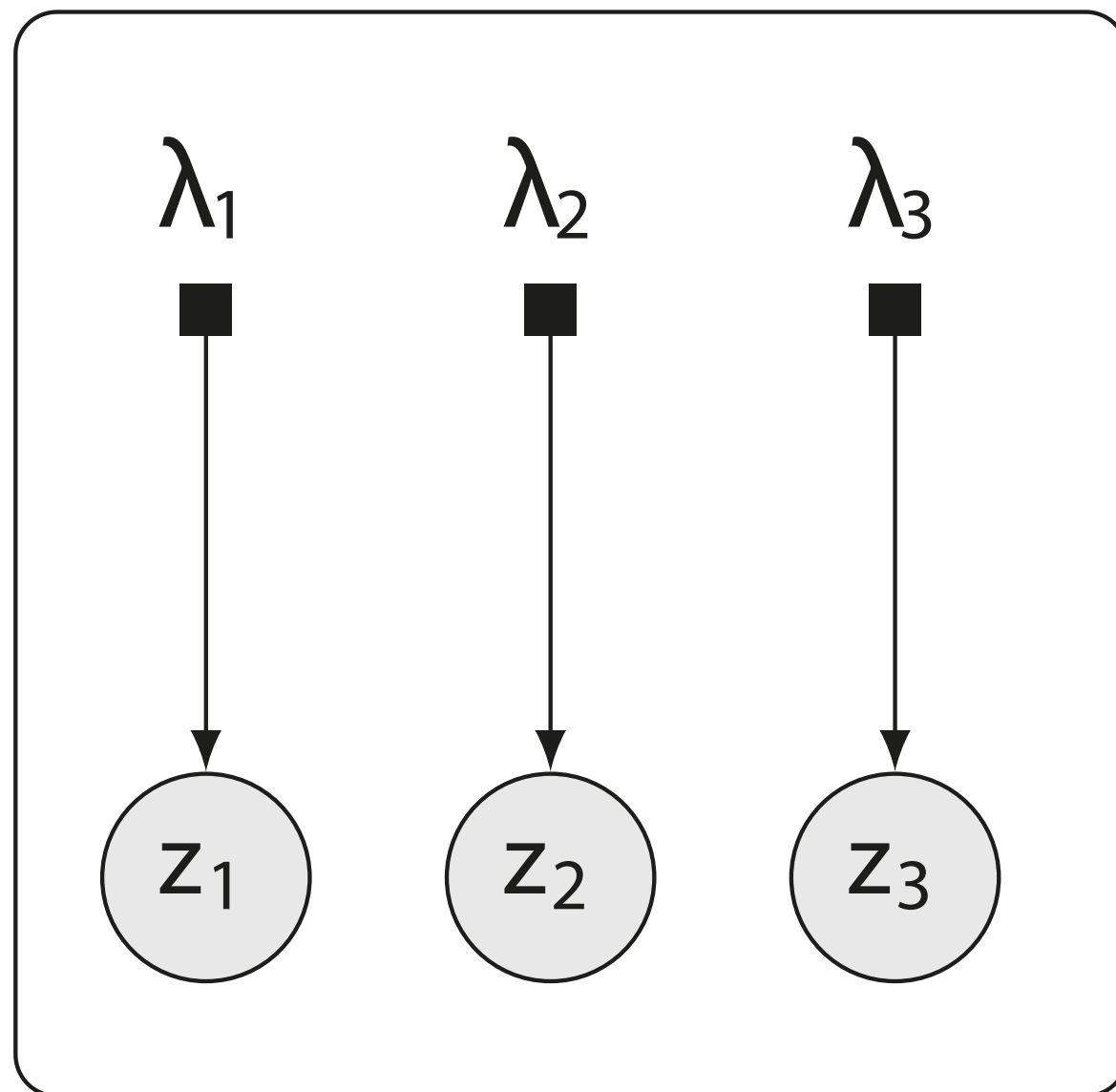
Variational Inference

$$\text{KL} (q(\mathbf{z}; \boldsymbol{\lambda}) \parallel p(\mathbf{z})) = \mathbb{E}_q[\log q(\mathbf{z}; \boldsymbol{\lambda})] - \mathbb{E}_q[-\beta \mathbf{H}(\mathbf{z})] + \log \mathcal{Z}$$



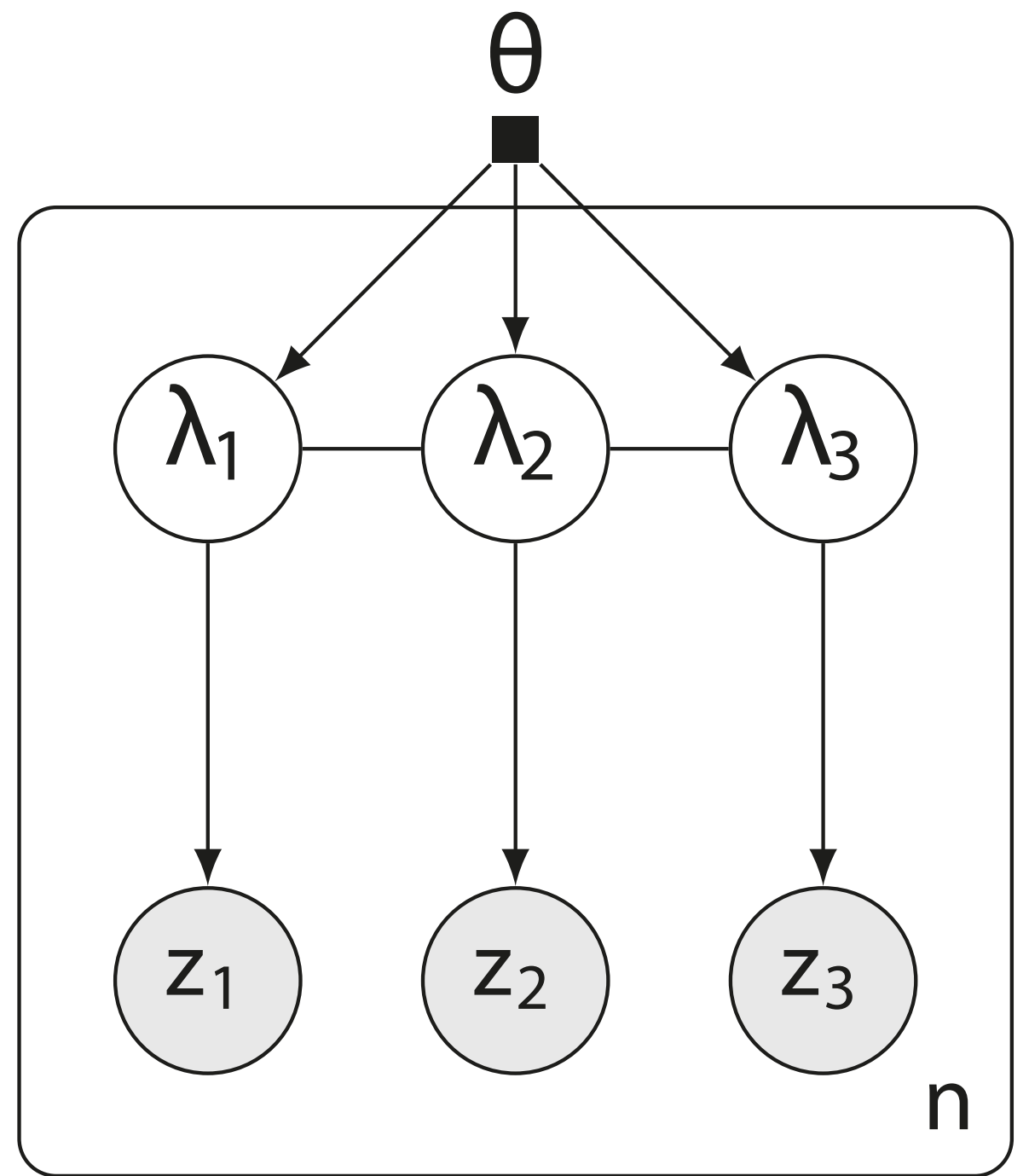
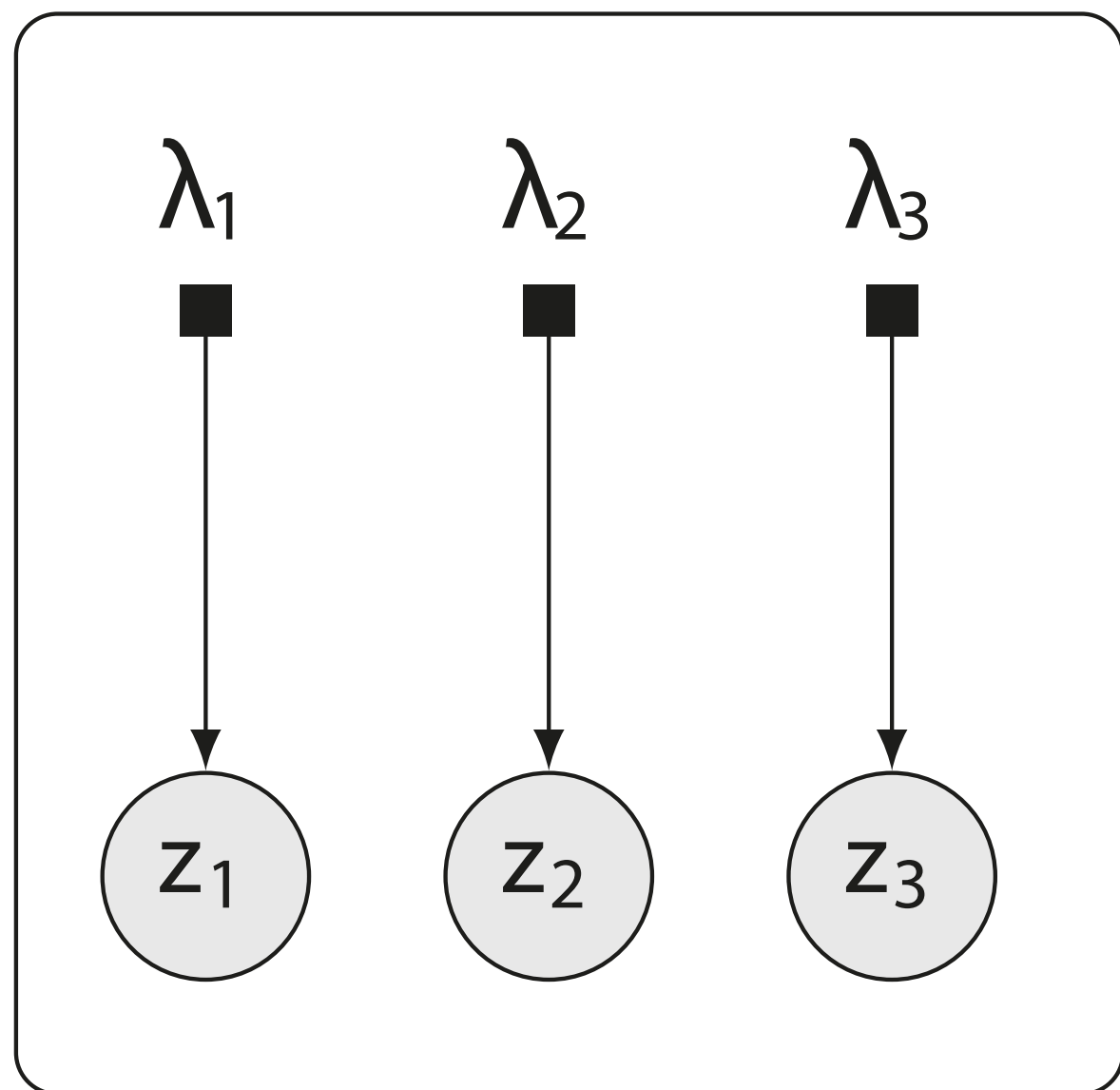
Mean Field Variational Family

$$q(\mathbf{z}; \boldsymbol{\lambda}) = \prod_{i=1}^d q(\mathbf{z}_i; \boldsymbol{\lambda}_i)$$

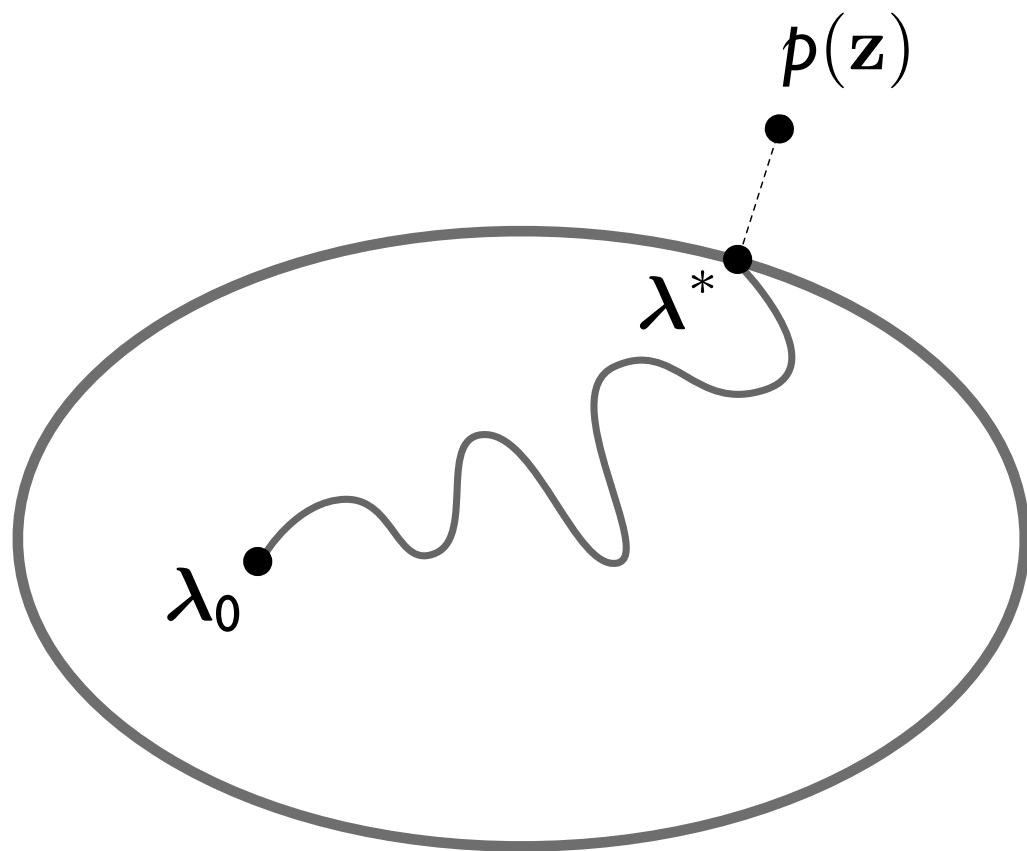


Hierarchical Variational Models

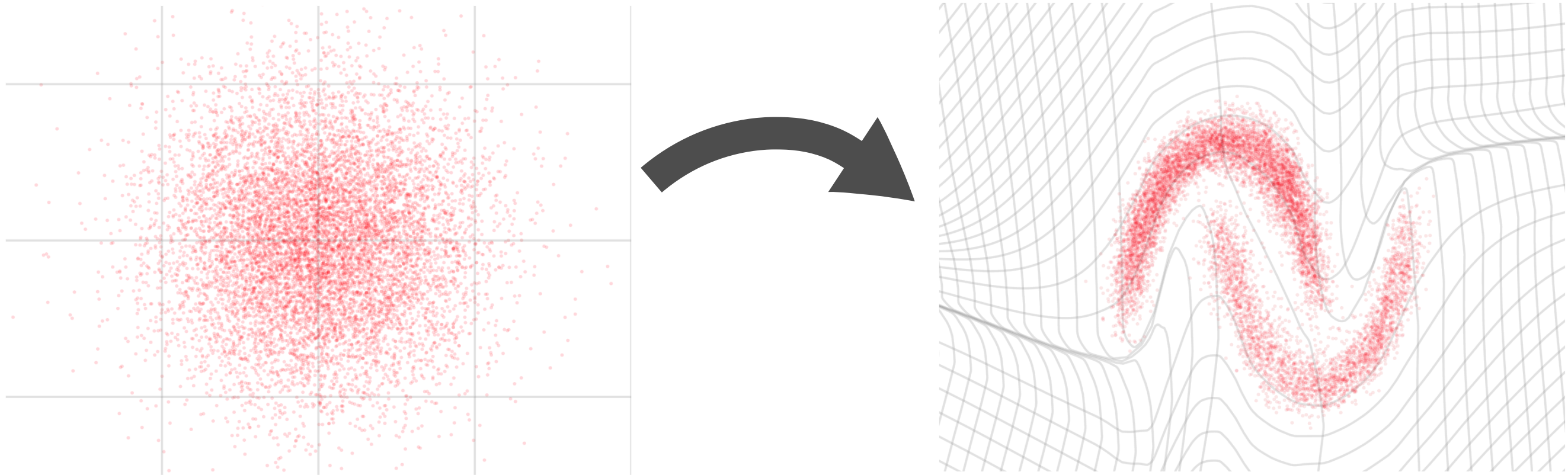
$$q(\mathbf{z}; \boldsymbol{\theta}) = \int \prod_{i=1}^d q(\mathbf{z}_i | \boldsymbol{\lambda}_i) q(\boldsymbol{\lambda}; \boldsymbol{\theta}) d\boldsymbol{\lambda}$$



Recursive Variational Posterior



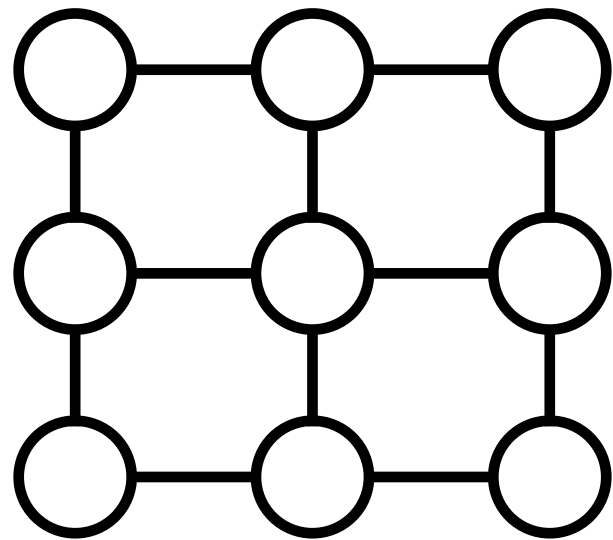
Generative Models: Flows



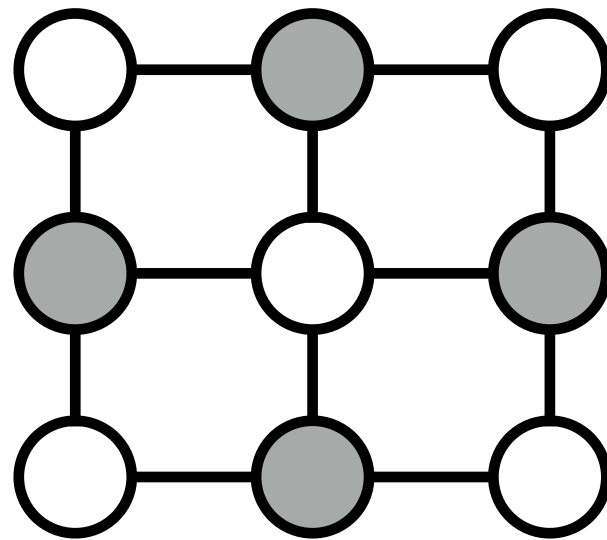
- Invertible transformations of random variables
- Cheap log density evaluation [Rezende+ 2015; Dinh+ 2016]
- In statistical physics: [Wu+ Phys. Rev. Lett. (2019)]
 - Variational autoregressive networks use PixelCNN [van den Oord+ 2016]

Building Hierarchical Variational Models

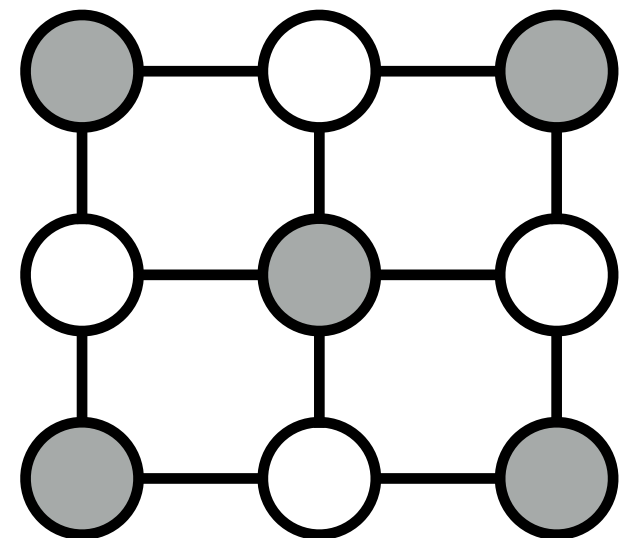
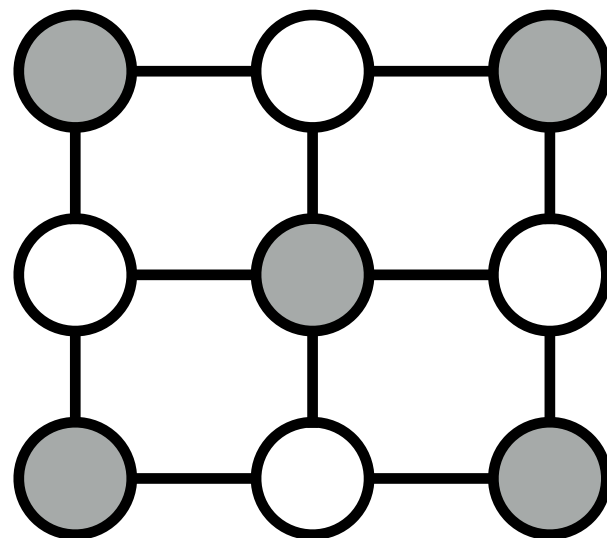
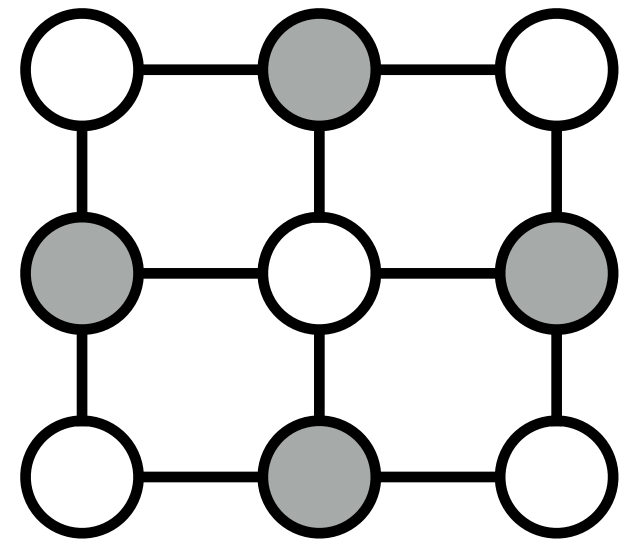
$$p(\mathbf{z})$$



$$q(\lambda; \theta)$$

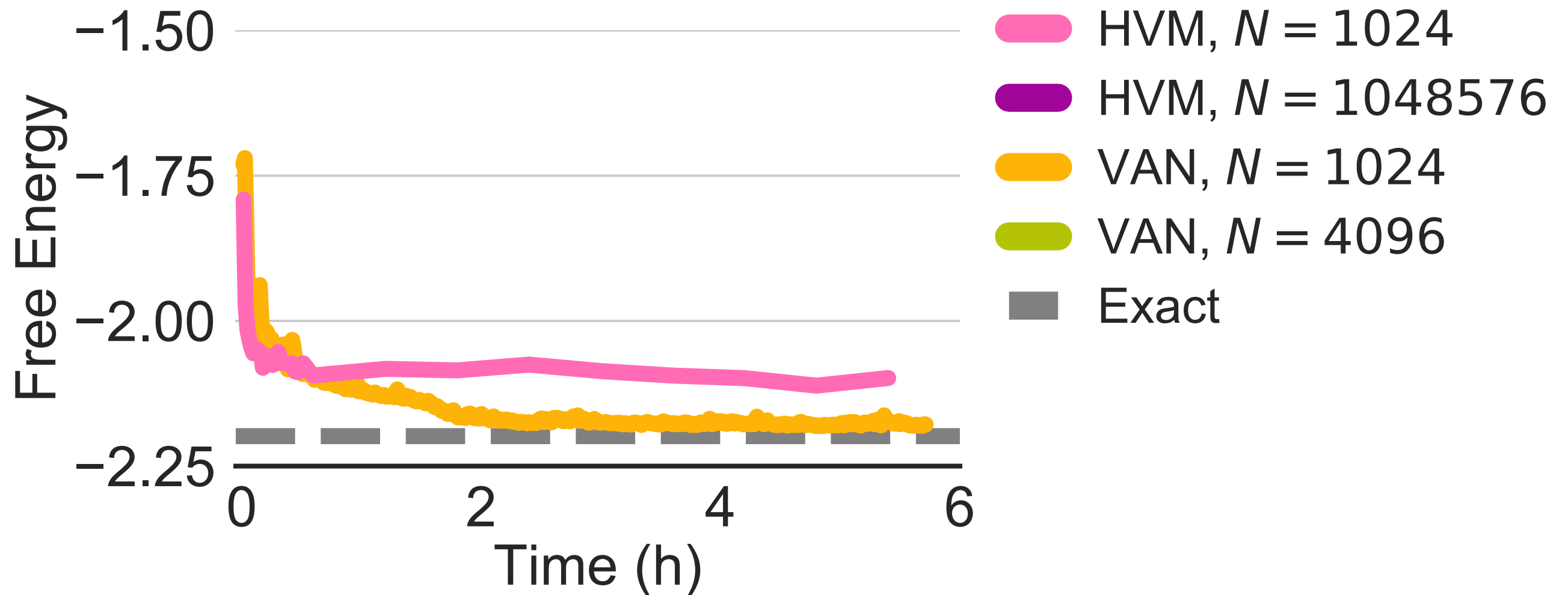


$$r(\lambda | \mathbf{z}; \phi)$$



Ising Model

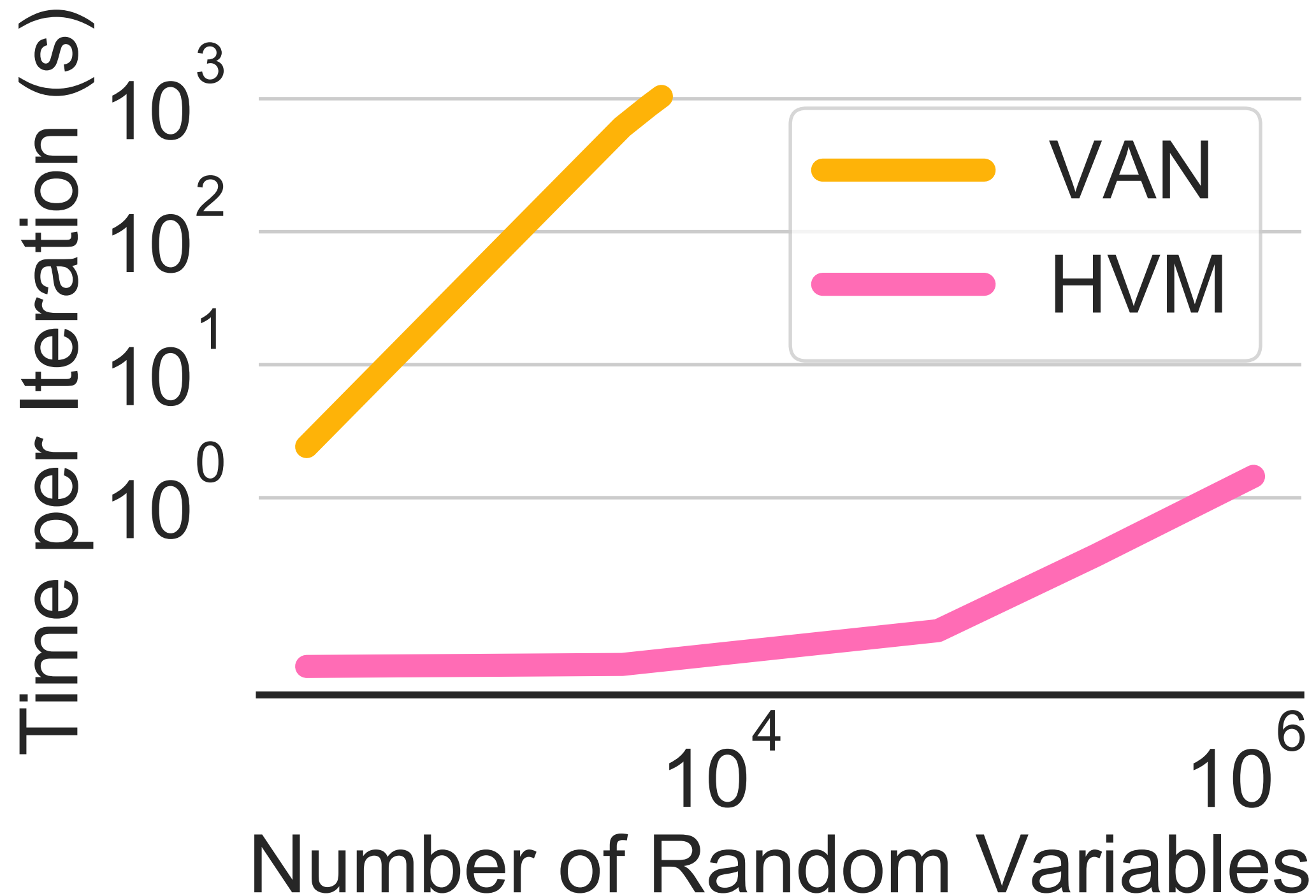
- HVM: 5400 parameters; VAN: 700k+
- Free energy: lower is better

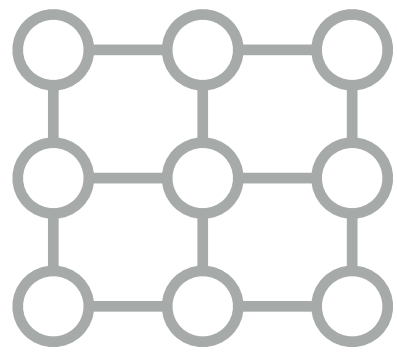


Sherrington-Kirkpatrick Model

- HVM: 5400 parameters; VAN: 700k+ parameters
- Free energy: lower is better

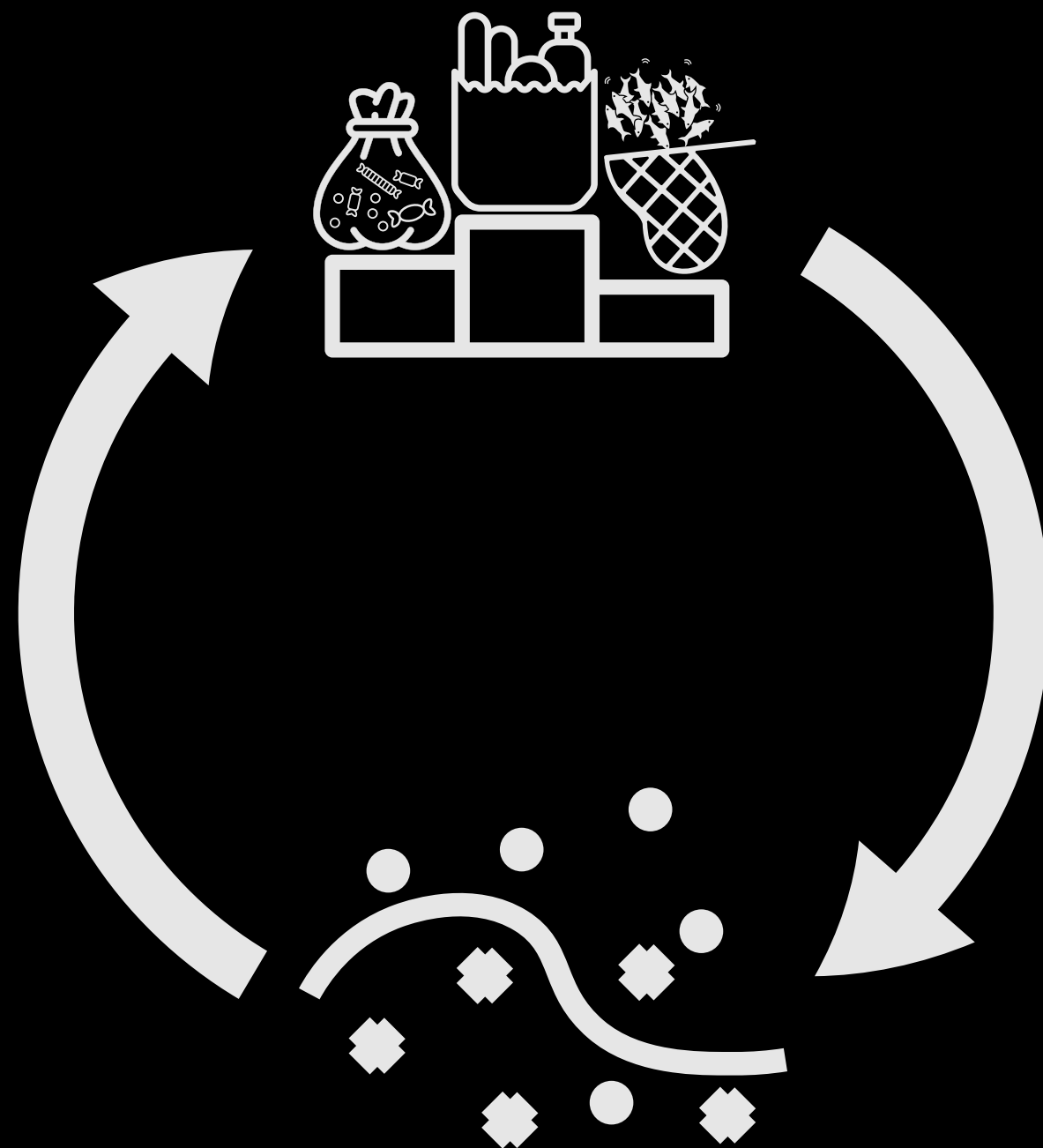
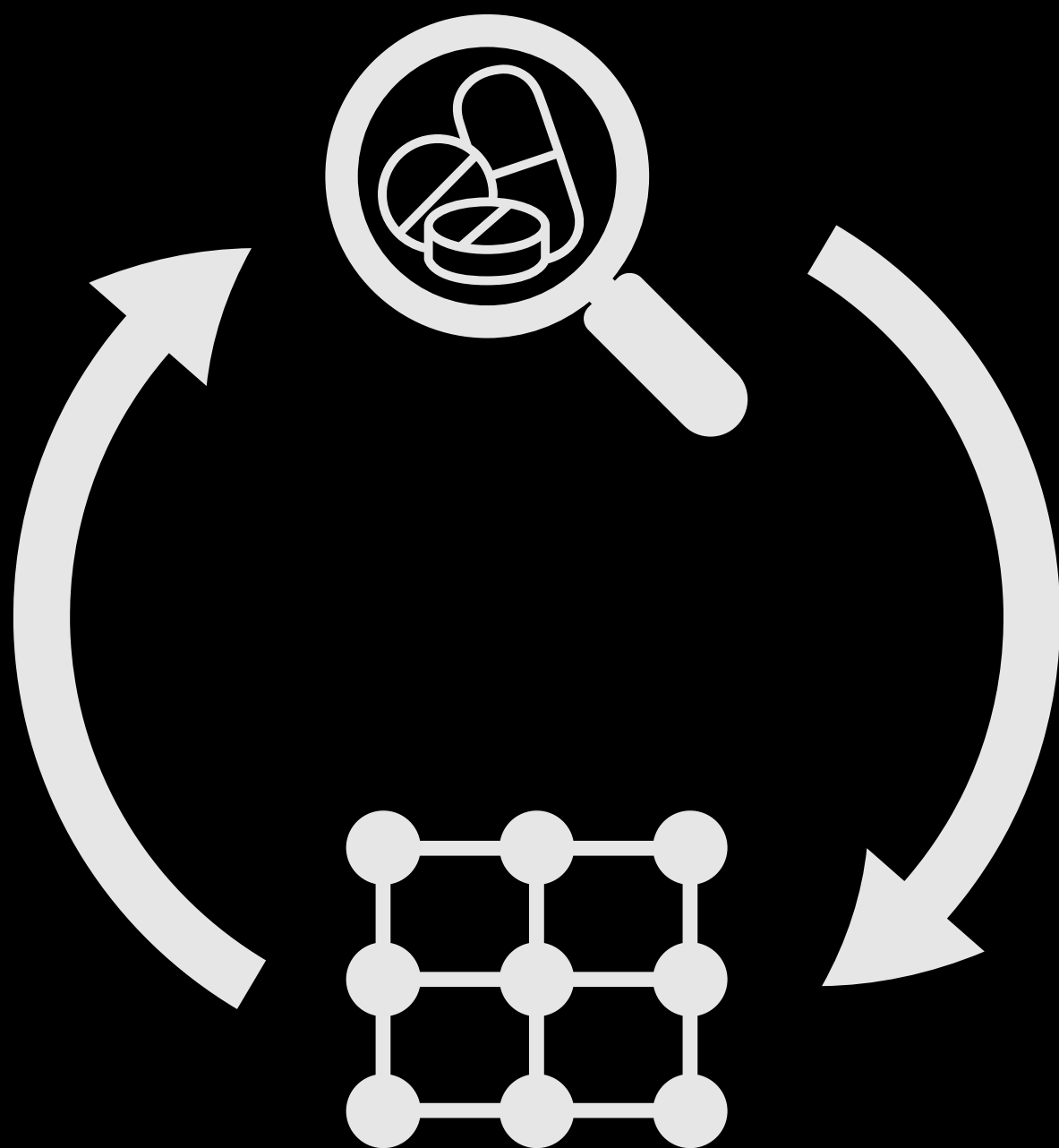
Scaling to Large Systems





Next Steps

- Use hierarchical variational models as proposal for Monte Carlo methods to reduce burn-in time
- Test on quantum systems for drug screening
- Rao-Blackwellization helps → increase accuracy with tighter lower bounds [Tucker+ 2018]
- Proximity variational inference can also increase accuracy [Altosaar+ 2018]









Emma Coats

@lawnrocket

Follow

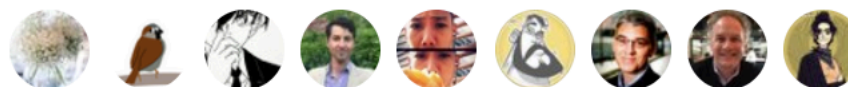


#4: Once upon a time there was _____. Every day, _____. One day _____. Because of that, _____. Because of that, _____. Until finally _____.

[#storybasics](#)

11:37 AM - 11 May 2012

34 Retweets 33 Likes



4

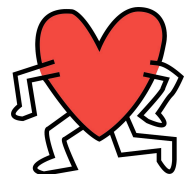


34



33

Once upon a time **there was matrix factorization.**

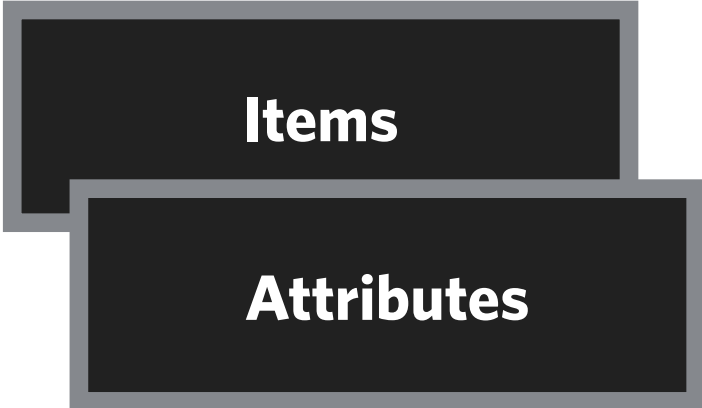
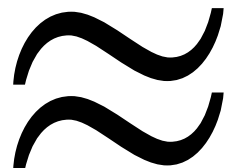


0	0	0	0	0	1	0
User-item interactions						
0	0	1	0	0	0	0

Emma Coats



0	0	0	0	1	0	0
Item attributes						
0	0	0	0	1	0	0



Once upon a time there was matrix factorization. Every day,
a new matrix factorization model appears.

Matrix factorization techniques for recommender systems

[PDF] datajobs.com

[Y Koren](#), [R Bell](#), [C Volinsky](#) - [Computer](#), 2009 - [computer.org](#)

[Findit@PUL](#)

AI

fa

re

te

☆

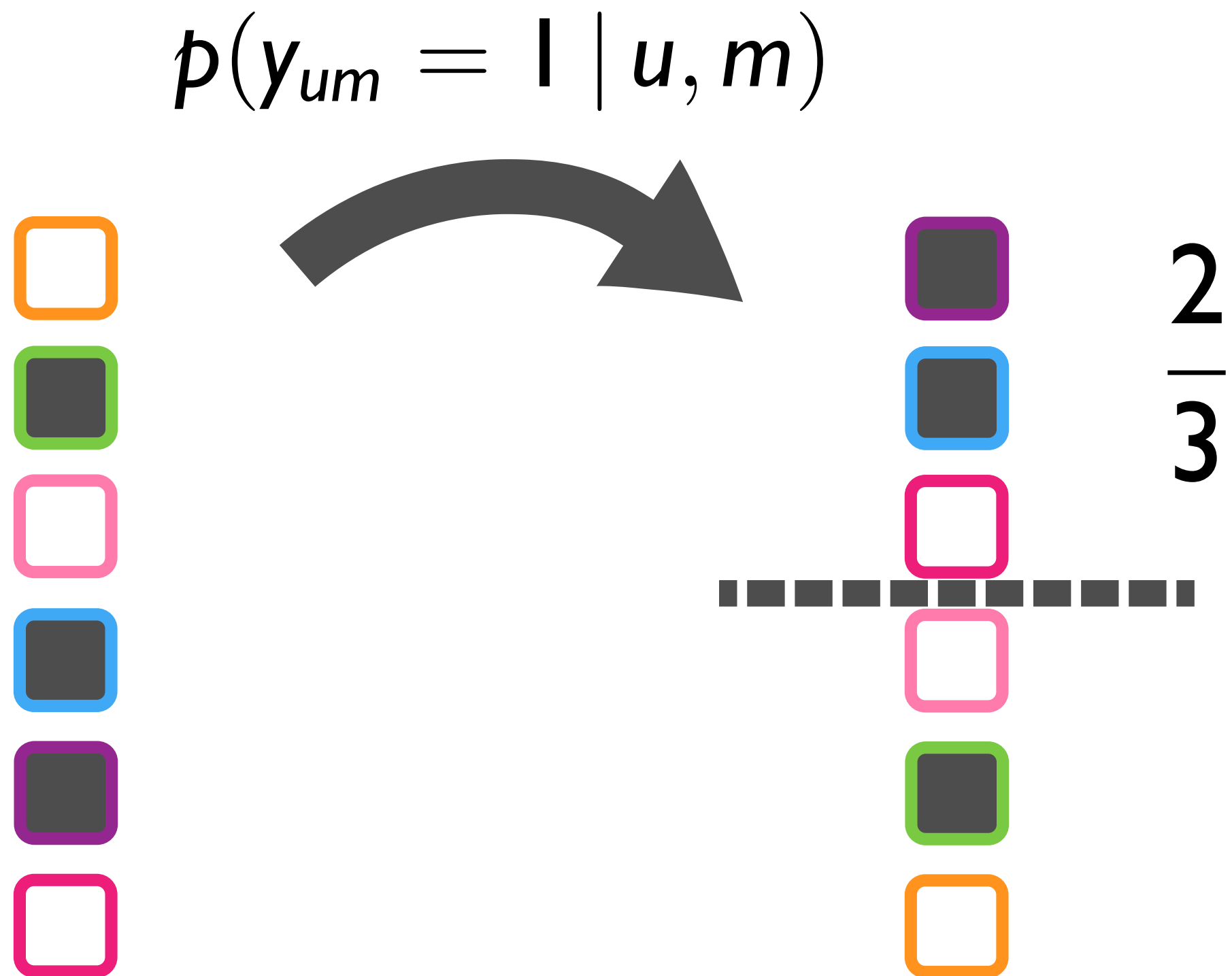
Factorization meets the item embedding: Regularizing matrix factorization with item-item co-occurrence

[D Liang](#), [J Altosaar](#), [L Charlin](#), [DM Blei](#) - ... of the 10th ACM conference on ..., 2016 - [dl.acm.org](#)

Matrix factorization (MF) models and their extensions are standard in modern recommender systems. MF models decompose the observed user-item interaction matrix into user and item latent factors. In this paper, we propose a co-factorization model, CoFactor, which jointly decomposes the user-item interaction matrix and the item-item co-occurrence matrix with shared item latent factors. For each pair of items, the co-occurrence matrix encodes the number of users that have consumed both items. CoFactor is inspired by the recent success ...

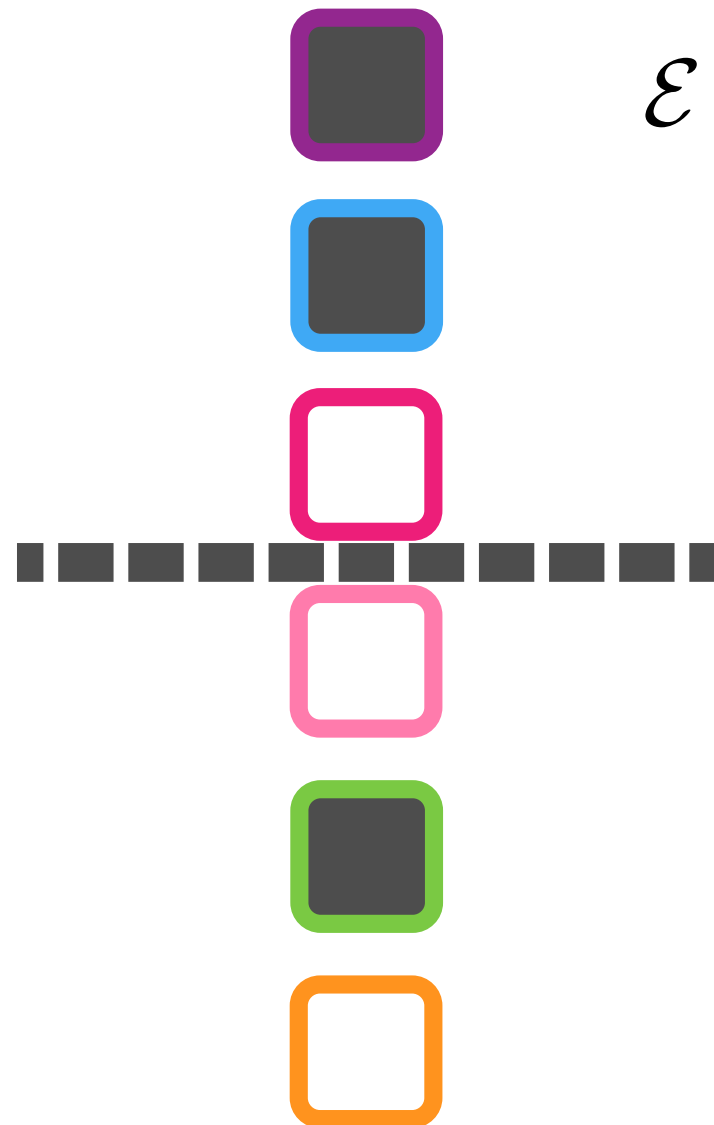
☆ [Cited by 137](#) [Related articles](#) [All 10 versions](#) [»»](#)

Once upon a time there was matrix factorization. Every day, a new matrix factorization model appears. One day, **recall** was used for evaluation.



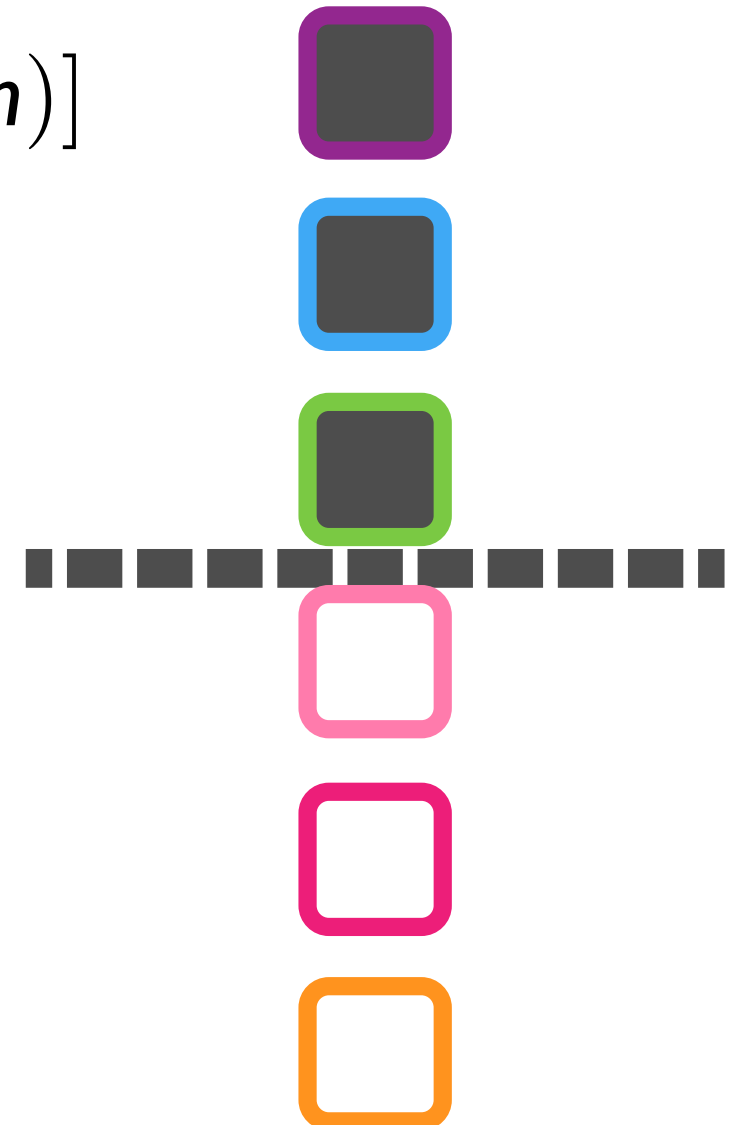
Once upon a time there was matrix factorization. Every day, a new matrix factorization model appears. One day recall was used for evaluation. Because of that, **zero worst-case error classifiers are optimal.**

$$\mathcal{E} > 0$$

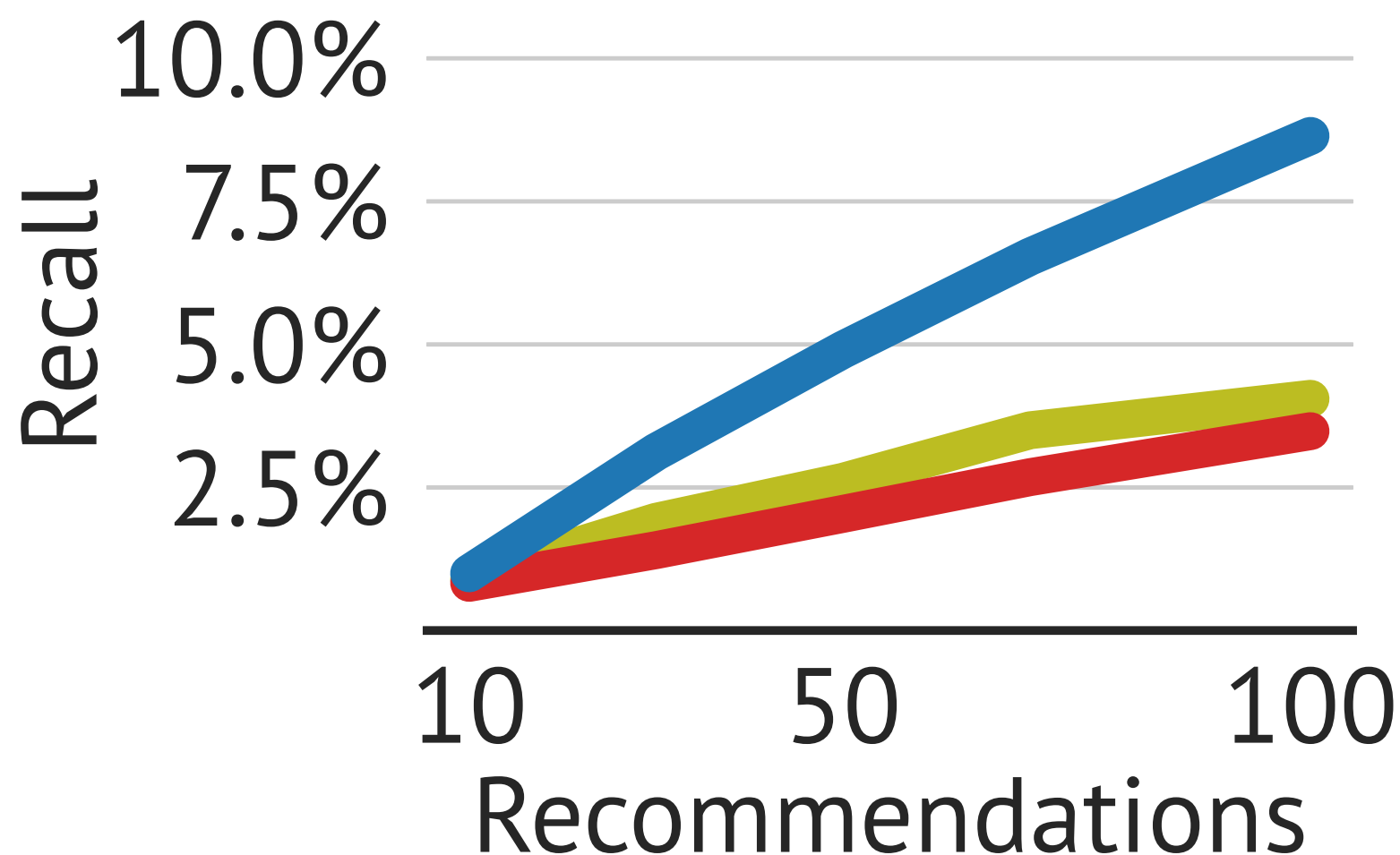


$$\mathcal{E} = \max_{(u,m) \in \mathcal{D}} \mathbb{I} [\hat{y}(u, m) \neq y(u, m)]$$

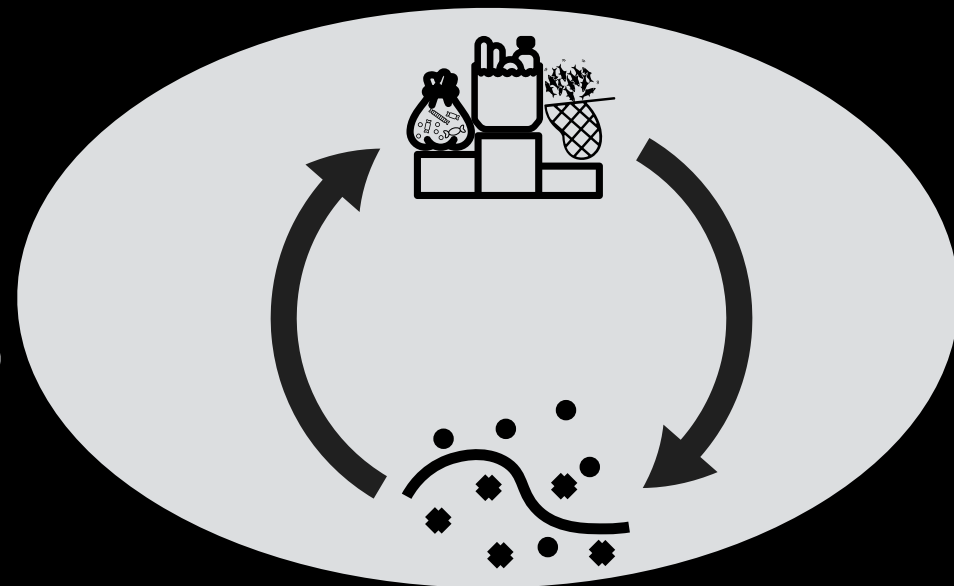
$$\mathcal{E} = 0$$

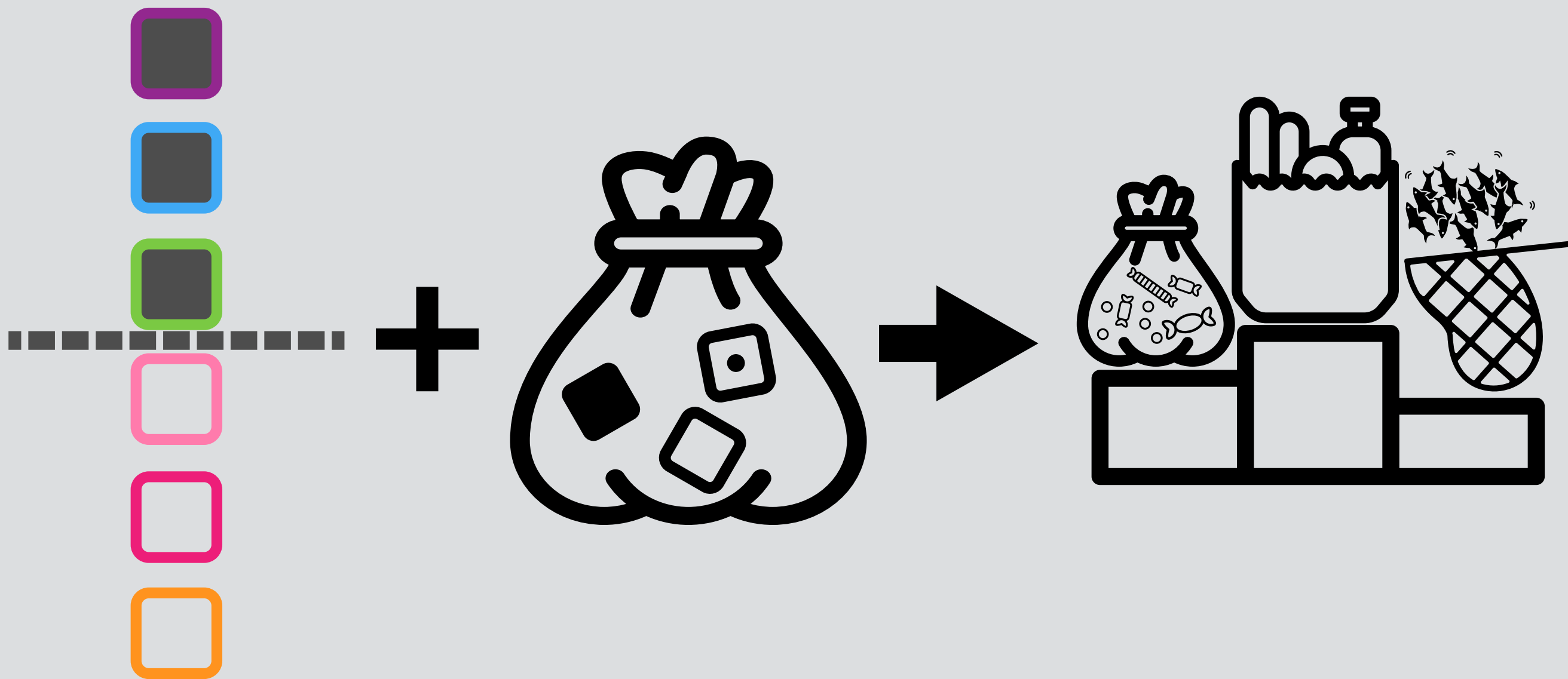


Once upon a time there was matrix factorization. Every day, a new matrix factorization model appears. One day recall was used for evaluation. Because of that, zero worst-case error classifiers are optimal. Because of that, **they outperform matrix factorization.**



Once upon a time there was matrix factorization. Every day, a new matrix factorization model appears. One day recall was used for evaluation. Because of that, zero worst-case error classifiers are optimal. Because of that, they outperform matrix factorization. Until finally **Jaan thought twice before dismissing classifiers as recommenders.**







Items (meals)	Attributes (foods)									Users				
	Pizza	Eggs	Taco	Salad	Avocado	Chicken	Sardines	Beer	Coffee	1	2	3	4	5
Breakfast pizza with coffee	✓	✓							✓		✓		✓	
Dinner pizza	✓						✓	✓		✓			✓	
Small salad				✓	✓		✓					✓		✓
Big salad		✓		✓	✓	✓			✓			✓	✓	
Taco			✓		✓	✓			✓					✓
Sardine taco			✓				✓			✓	✓		✓	



RankFromSets

$$p(y_{um} = l \mid u, m) = \sigma(f(u, x_m))$$

- f is a neural network, invariant to permutation of x_m
- This choice of model can maximize recall [Altosaar+ 2020]



RankFromSets

$$p(y_{um} = l \mid u, m) = \sigma(f(u, x_m))$$

$$f(u, x_m) =$$



RankFromSets

$$p(y_{um} = l \mid u, m) = \sigma(f(u, x_m))$$

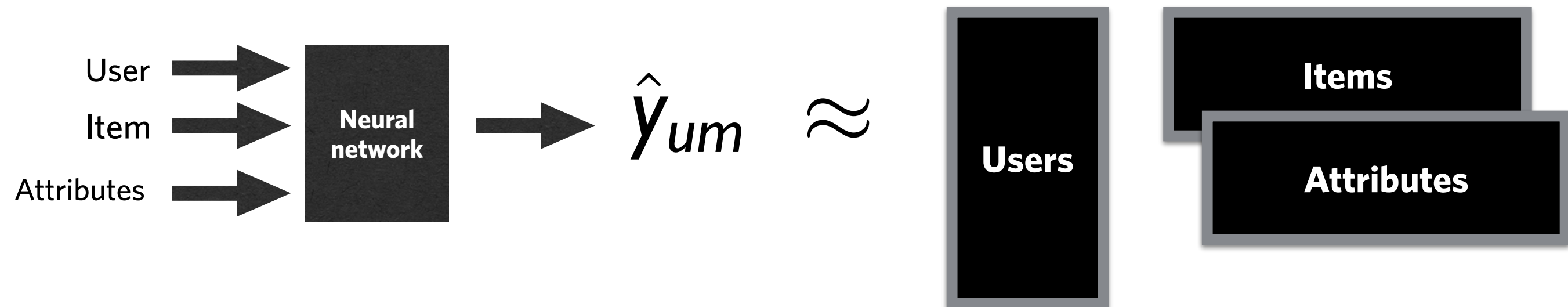


Universal approximation: RankFromSets can approximate any order-invariant model.

$$f(u, x_m)$$

Examples of models invariant to permutation of item attributes:

- Matrix factorization
- Permutation-marginalized recurrent neural networks



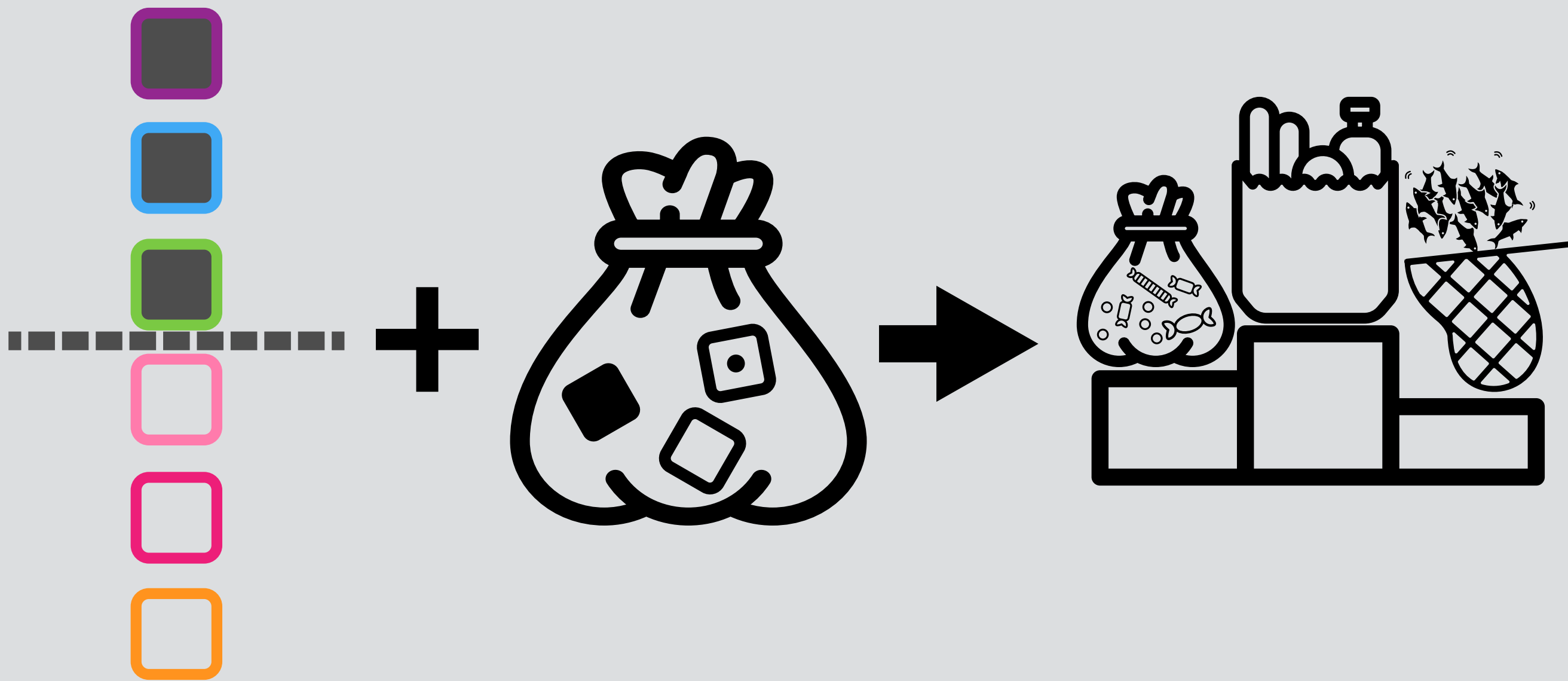
[Zaheer+ 2017; Altosaar+ 2020]



RankFromSets Objective

$$\mathcal{L}(\gamma, \lambda_u) = \mathbb{E}_u \left[\mathbb{E}_{m \sim \mathcal{D}_u \mid y_{um}=1} [\log p(y_{um} = 1 \mid \mathbf{x}_m; \gamma)] \right. \\ \left. + \lambda_u \mathbb{E}_{k \sim \mathcal{D}_u \mid y_{uk}=0} [\log p(y_{uk} = 0 \mid \mathbf{x}_k; \gamma)] \right]$$

- Parameters γ : user embeddings, attribute embeddings, item embeddings, weights and biases
- λ_u balances negative examples for every user

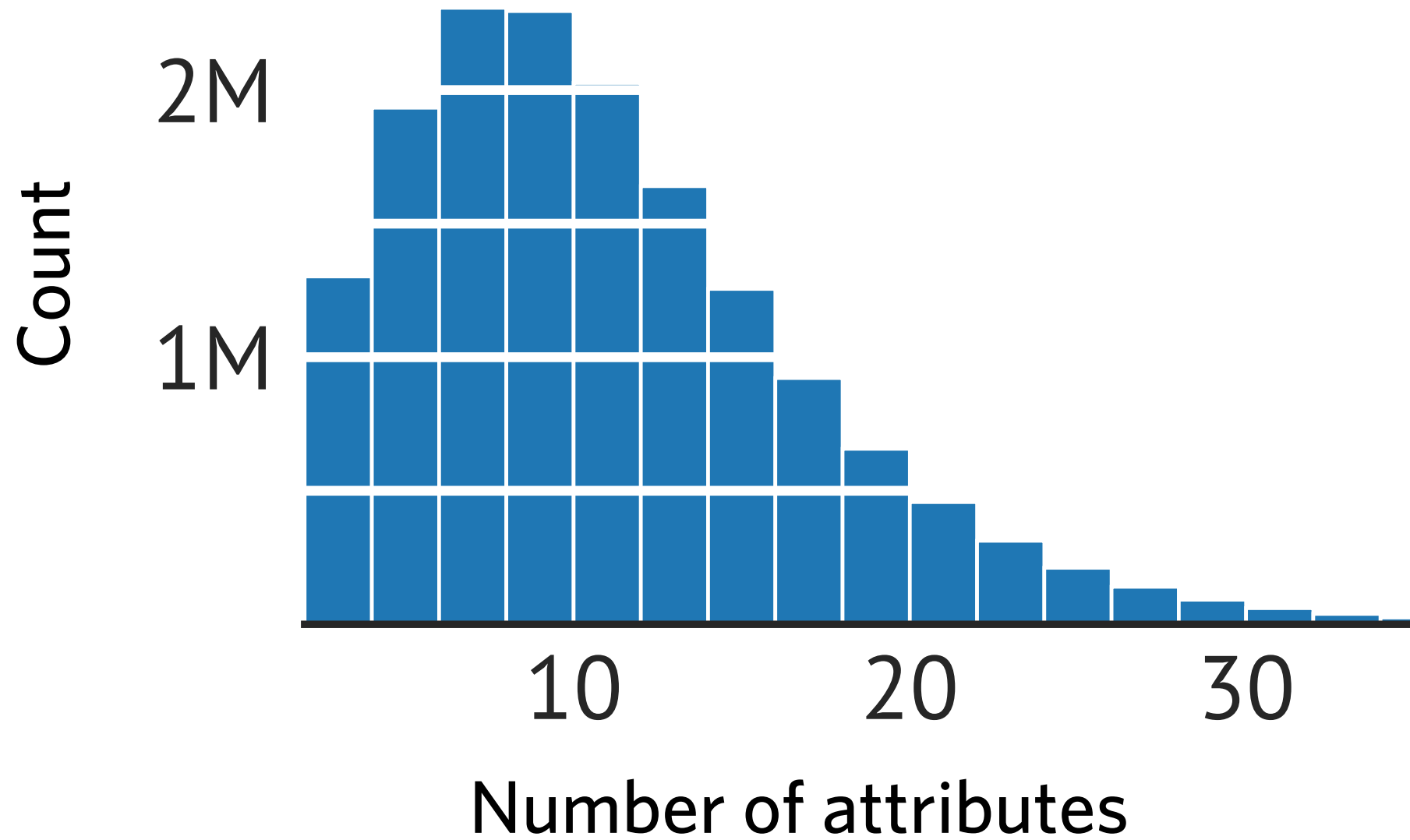








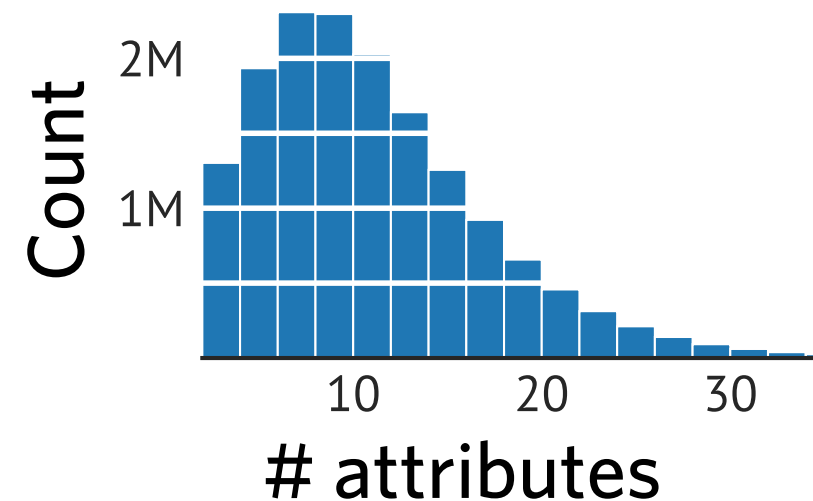
Lose It!





One year of data from 55,000 users of the Lose It! food tracking app

- 15 million meals
- Item attributes: 9,963 food words



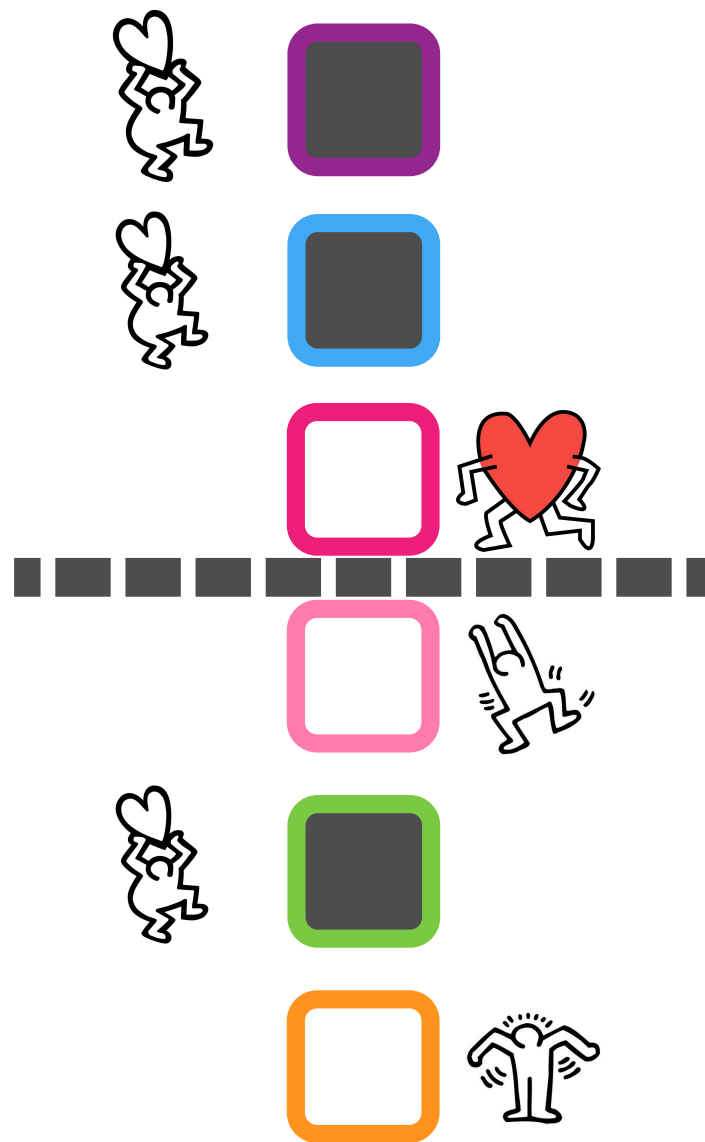
Example meals:

- *TWO SCOOPS OF RAISIN BRAN CEREAL, ORGANIC MOROCCAN GREEN TEA, ALMOND MILK, LIGHT HONEY, TAP WATER, LARGE BANANA, LARGE STRAWBERRIES*
- *BOSTON ROAST PORK, MACKEREL, ARTICHOKE HEARTS, SPINACH, PIMIENTO-STUFFED MANZANILLA OLIVES, CARROTS, MUSHROOMS, PEPPERCORN RANCH DRESSING*



Scalable Evaluation

For every held out item, sample negative labels from other users.





Lose It!

-  RFS f : residual
-  RFS f : inner product
-  CTPF
-  Word embedding
-  StarSpace
-  LSTM
-  Random

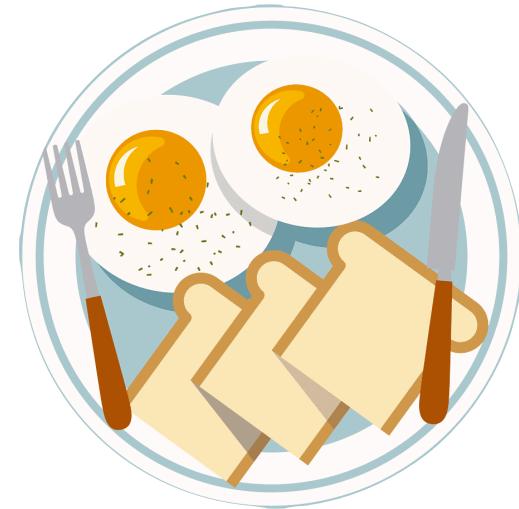
- CTPF: Collaborative Topic Poisson Factorization [Gopalan+ 2014]
- StarSpace [Wu+ 2018]
- LSTM [Bansal+ 2016]
- Word embedding [Bojanowski+ 2016]

Query



Two scoops of Raisin Bran cereal,
organic Moroccan green tea,
almond milk, light honey, tap water,
large banana, large strawberries

Nearest Neighbor



Vita Bee bread, salted butter, fresh
medium tomatoes, large fried
whole egg, small banana

Query

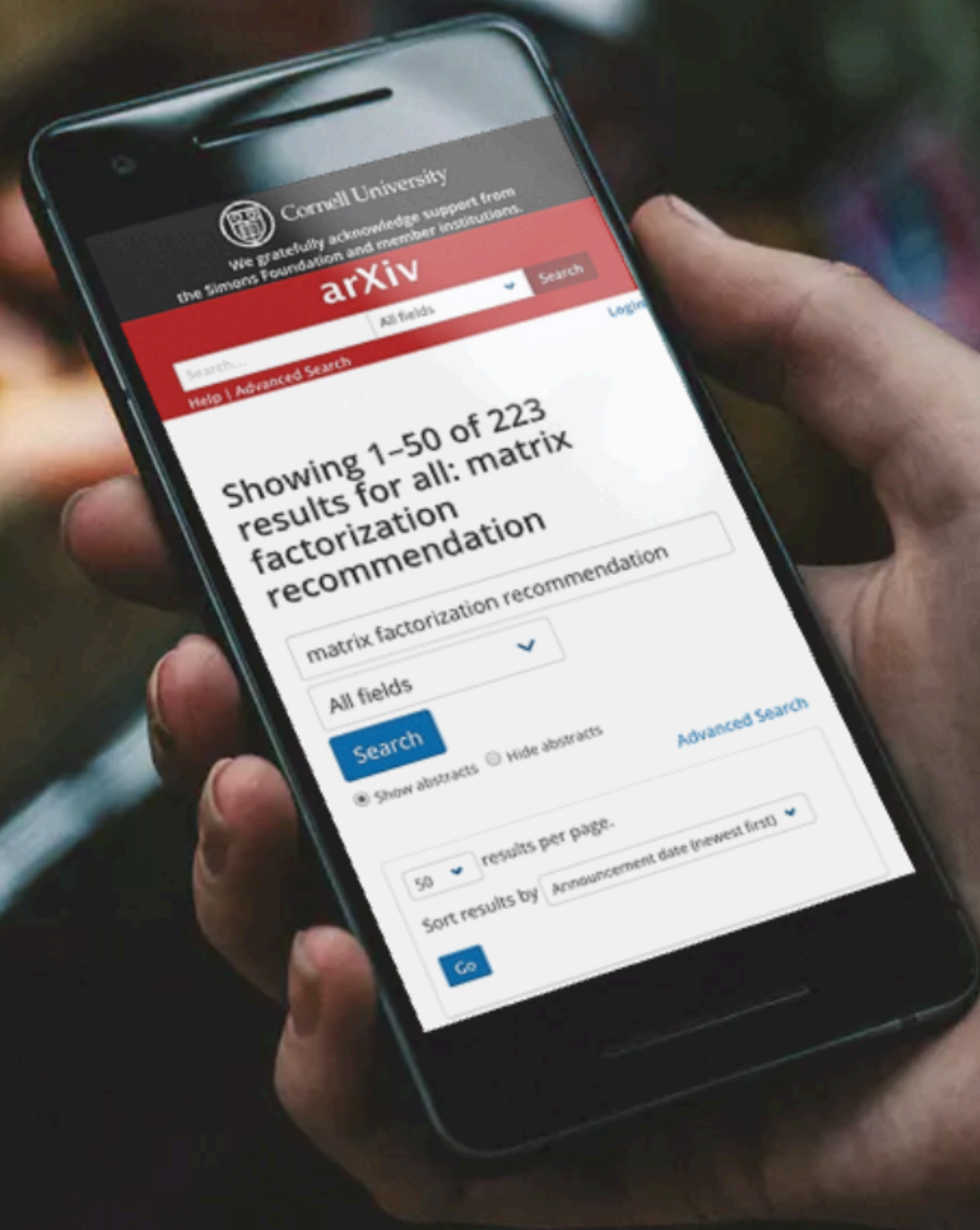


Iceberg lettuce, cantaloupe cubes, diced honeydew melon, cherry tomatoes, olives, dry-cooked unsalted hulled sunflower seed kernels, chopped hard-boiled egg, cucumbers, dried cranberries, fat-free ranch dressing

Nearest Neighbor



Green leaf lettuce, chopped sweet red bell peppers, crumbled feta cheese, large hard-boiled egg, chopped cucumber, oil-roasted salted sunflower seeds, sliced radishes, sliced strawberries, pitted Calamata olives, fat-free balsamic vinegar



Cornell University

We gratefully acknowledge support from the Simons Foundation and member institutions.

arXiv

All fields

Search

Search...
Help | Advanced Search

Login

Showing 1-50 of 223 results for all: matrix factorization recommendation

matrix factorization recommendation

All fields

Search

☒ Show abstracts ☐ Hide abstracts

Advanced Search

50 results per page.

Sort results by

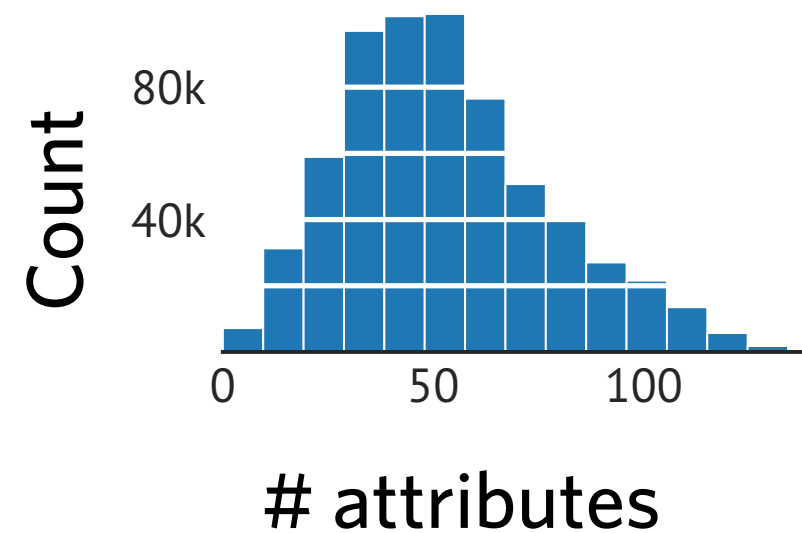
Announcement date (newest first)

Go



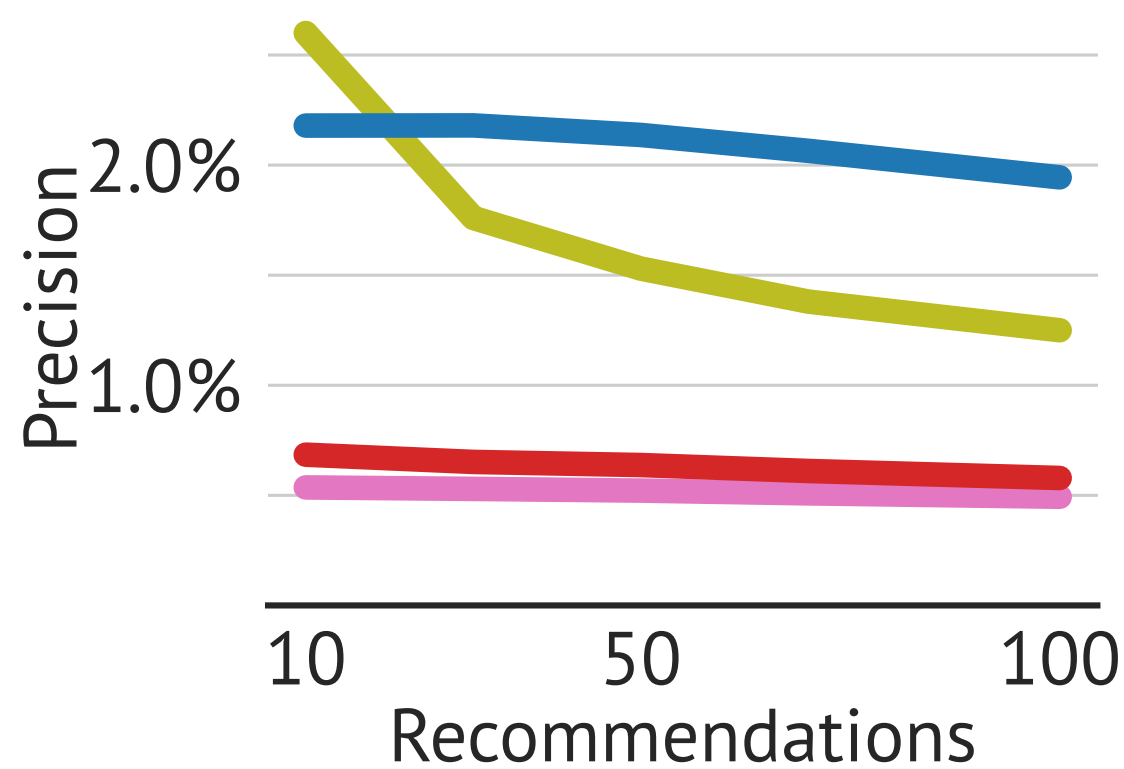
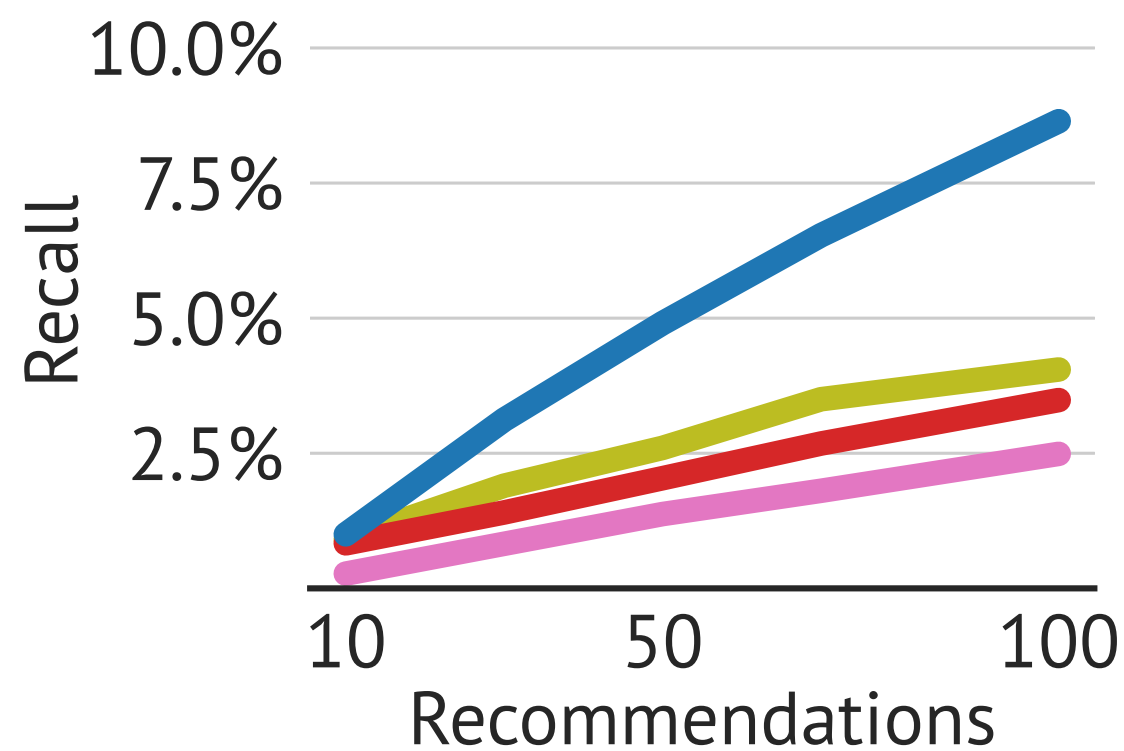
One year of data, 64,978 users

- 636,622 papers
- Vocabulary of 14,000 in abstracts

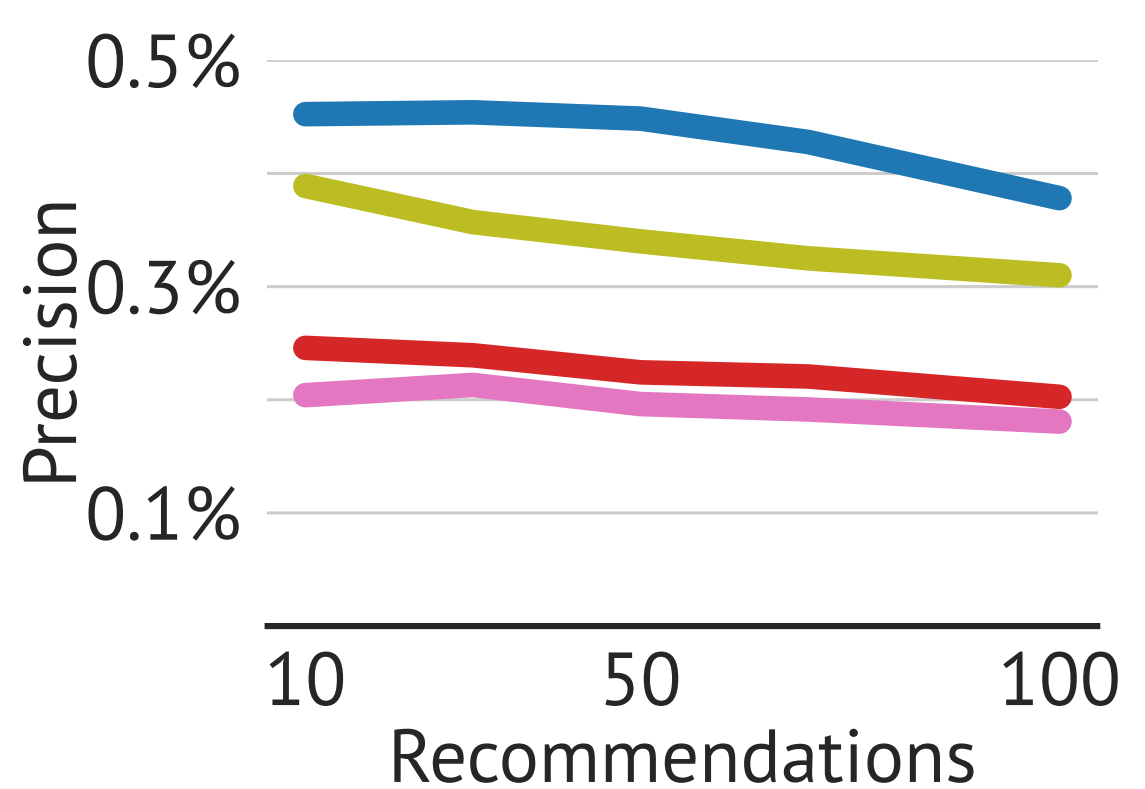
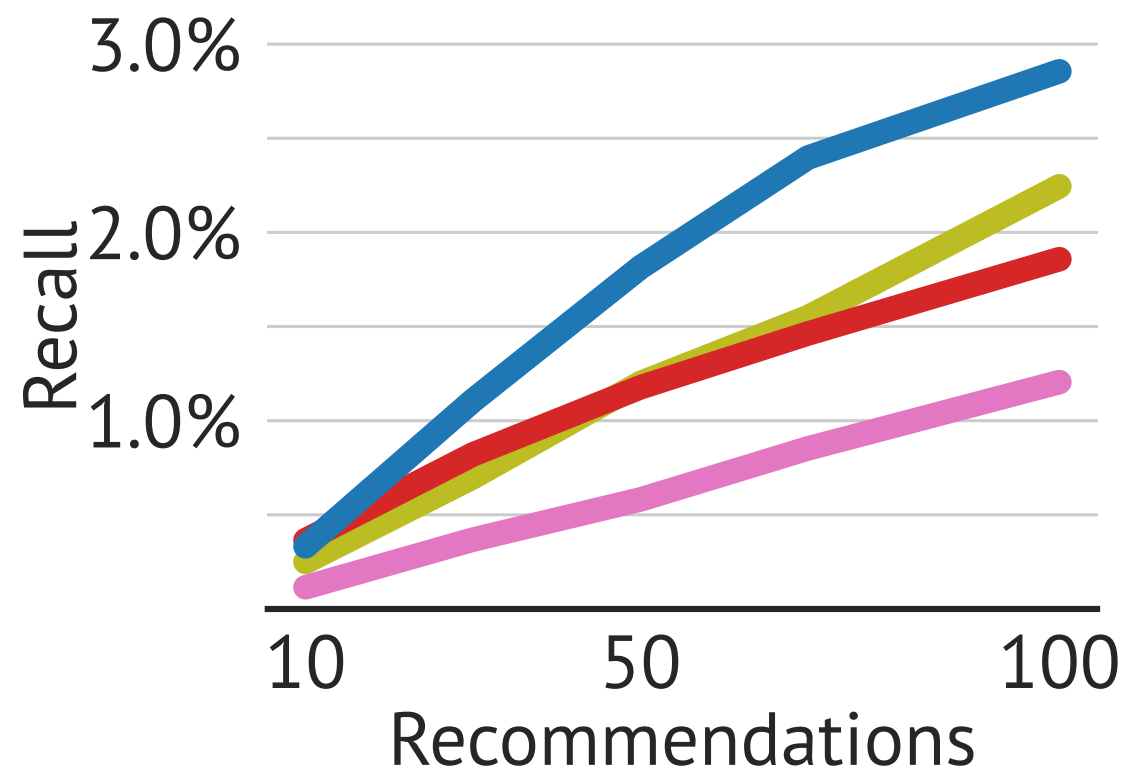


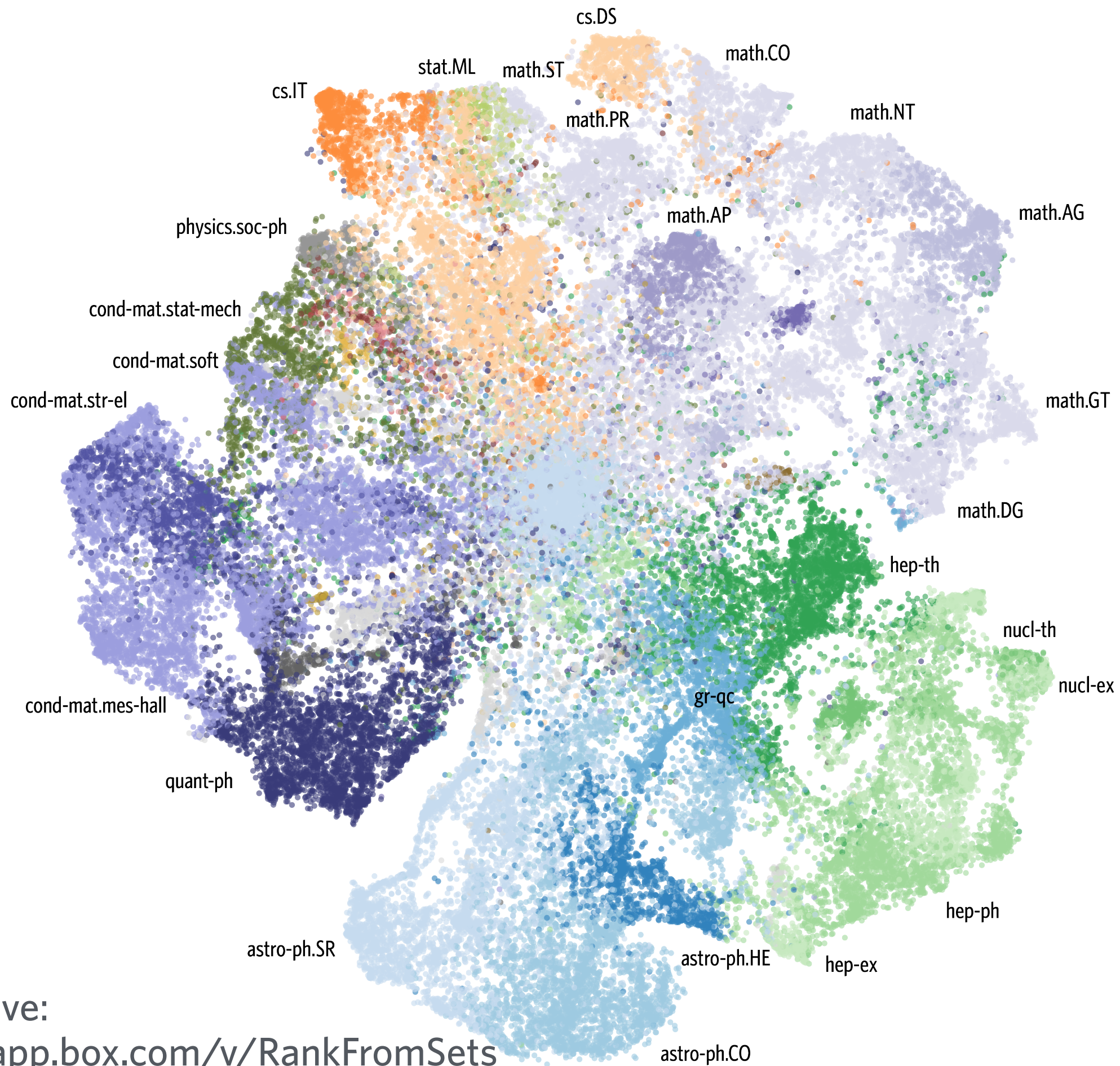
- RankFromSets
- CTPF
- Word embedding
- StarSpace

In-matrix

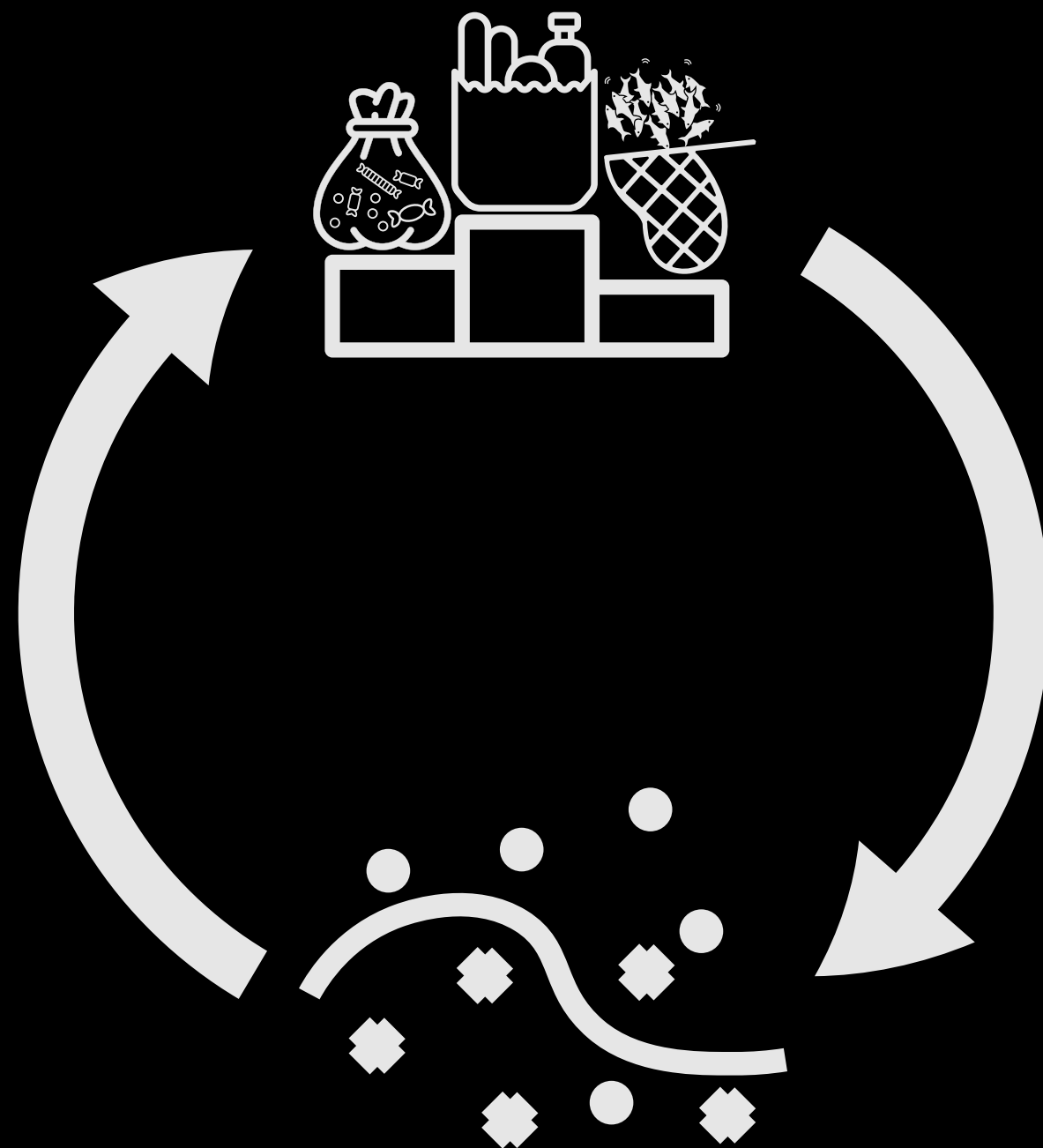
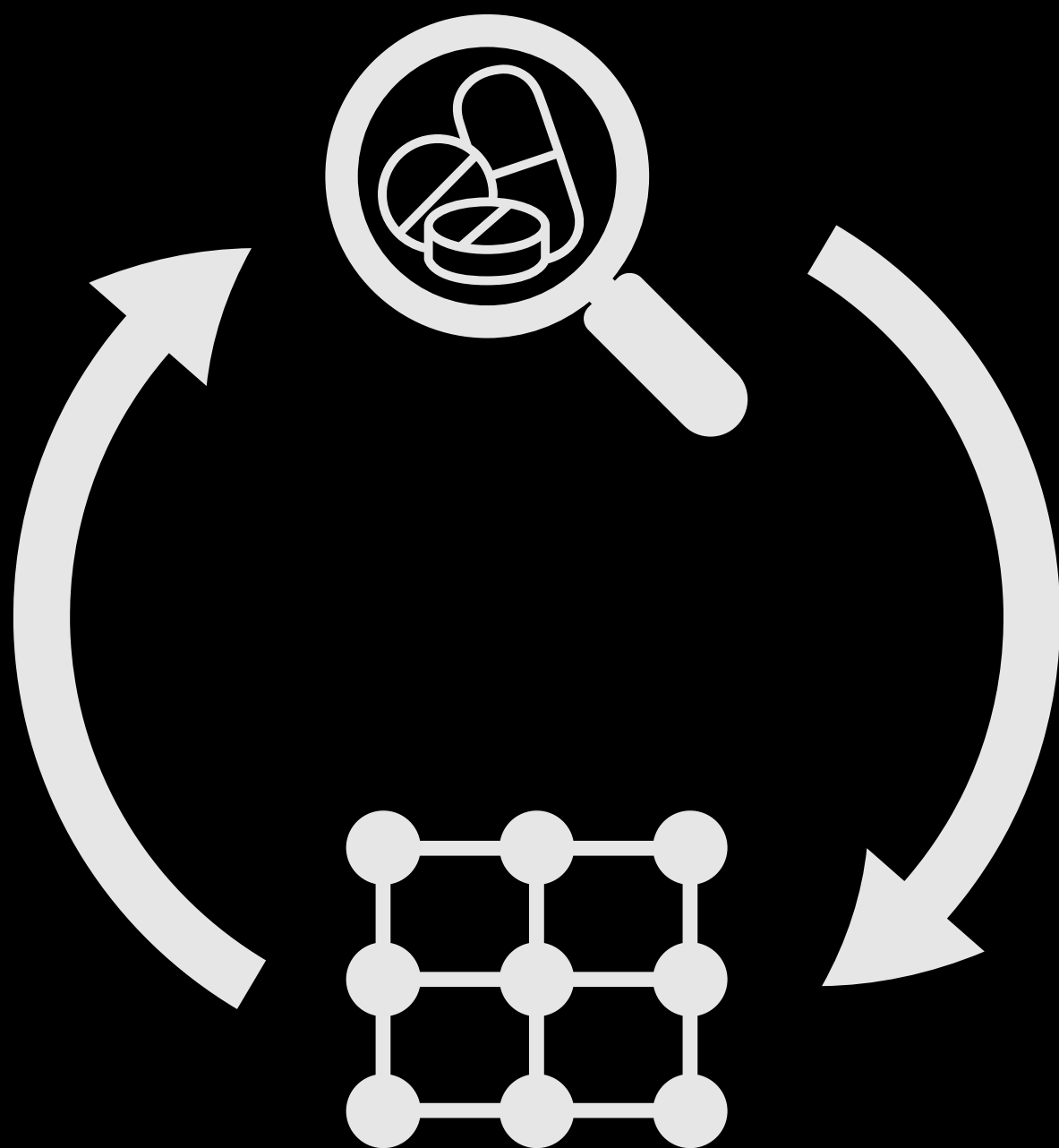


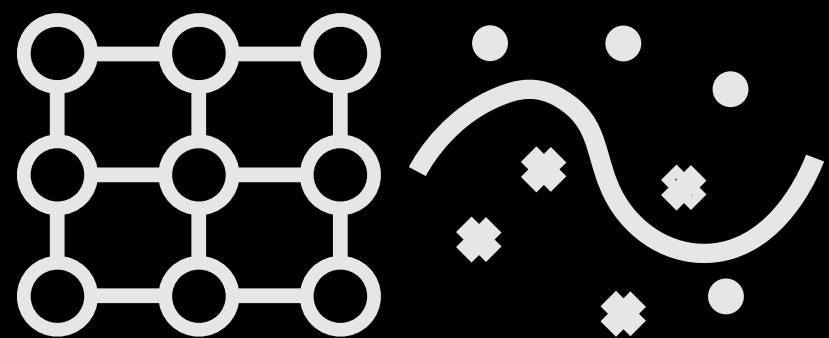
Out-matrix





Interactive:
<https://app.box.com/v/RankFromSets>





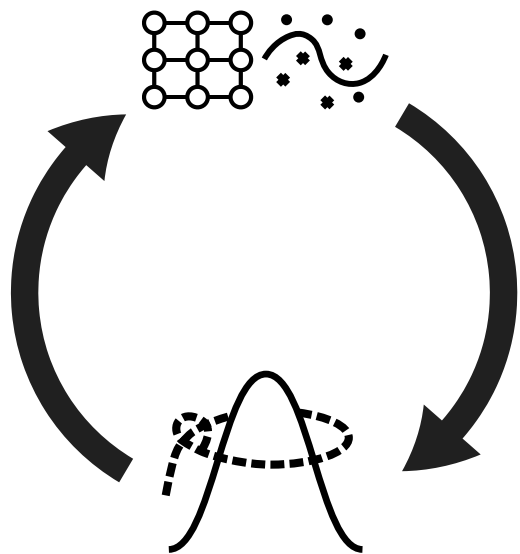
Bernoulli Factor Model

$$\mathbf{z}_{ik} \sim \text{Bernoulli}(\pi)$$

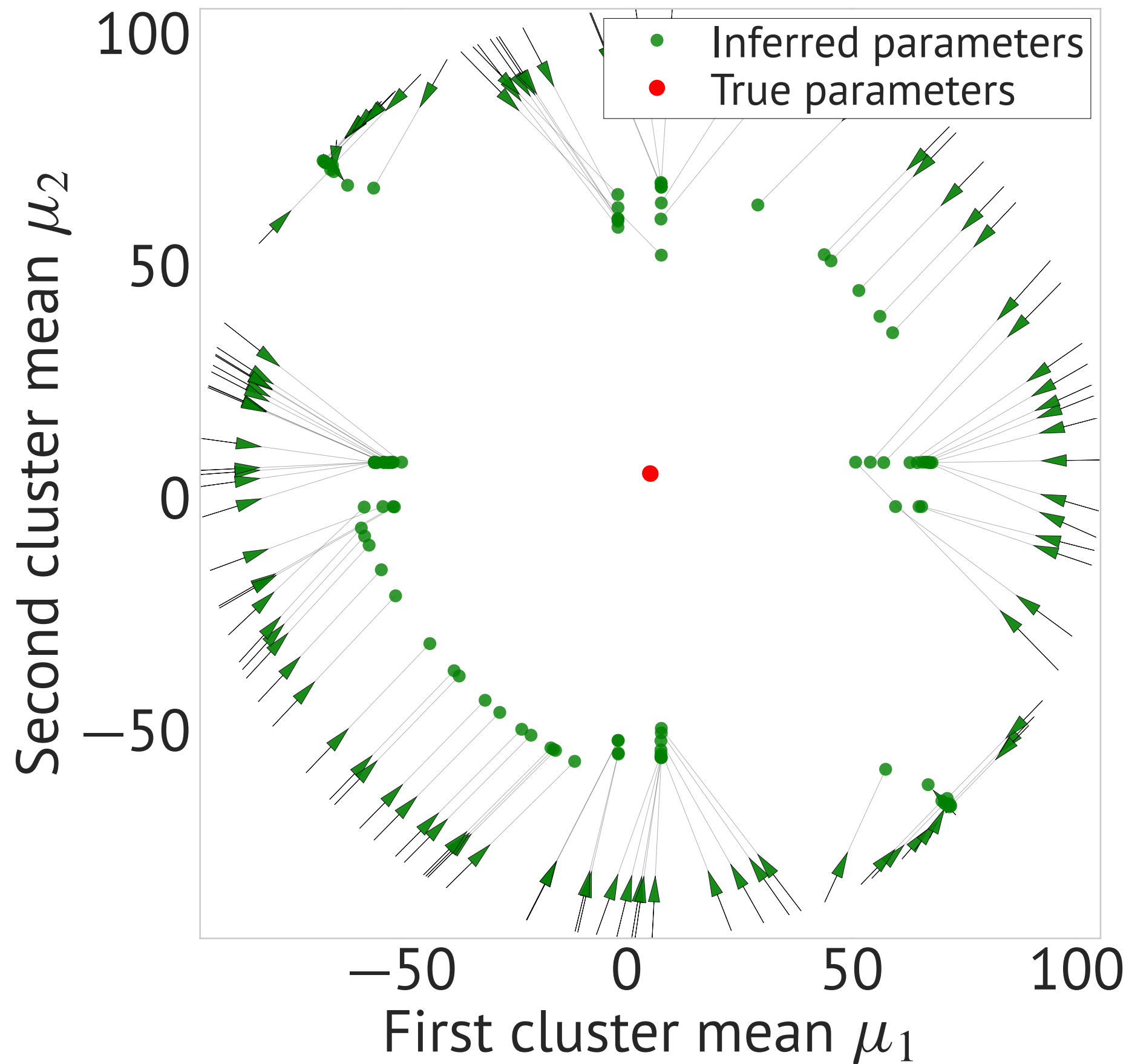
$$\mathbf{x}_i \sim \text{Normal}(\mathbf{z}_i^\top \boldsymbol{\mu}, \sigma^2 = \mathbf{I})$$



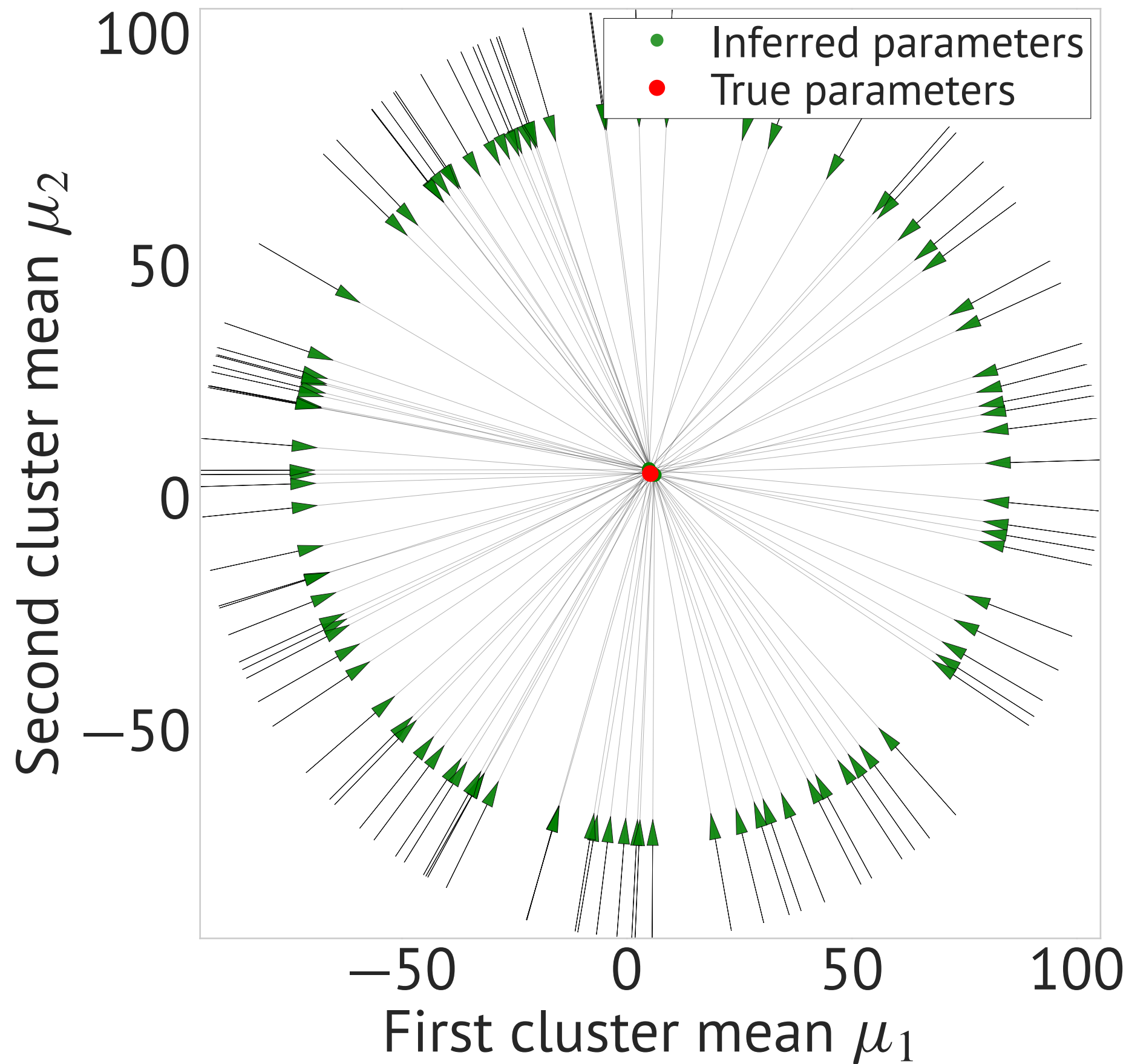
Goal: infer cluster mean parameter



Variational Inference

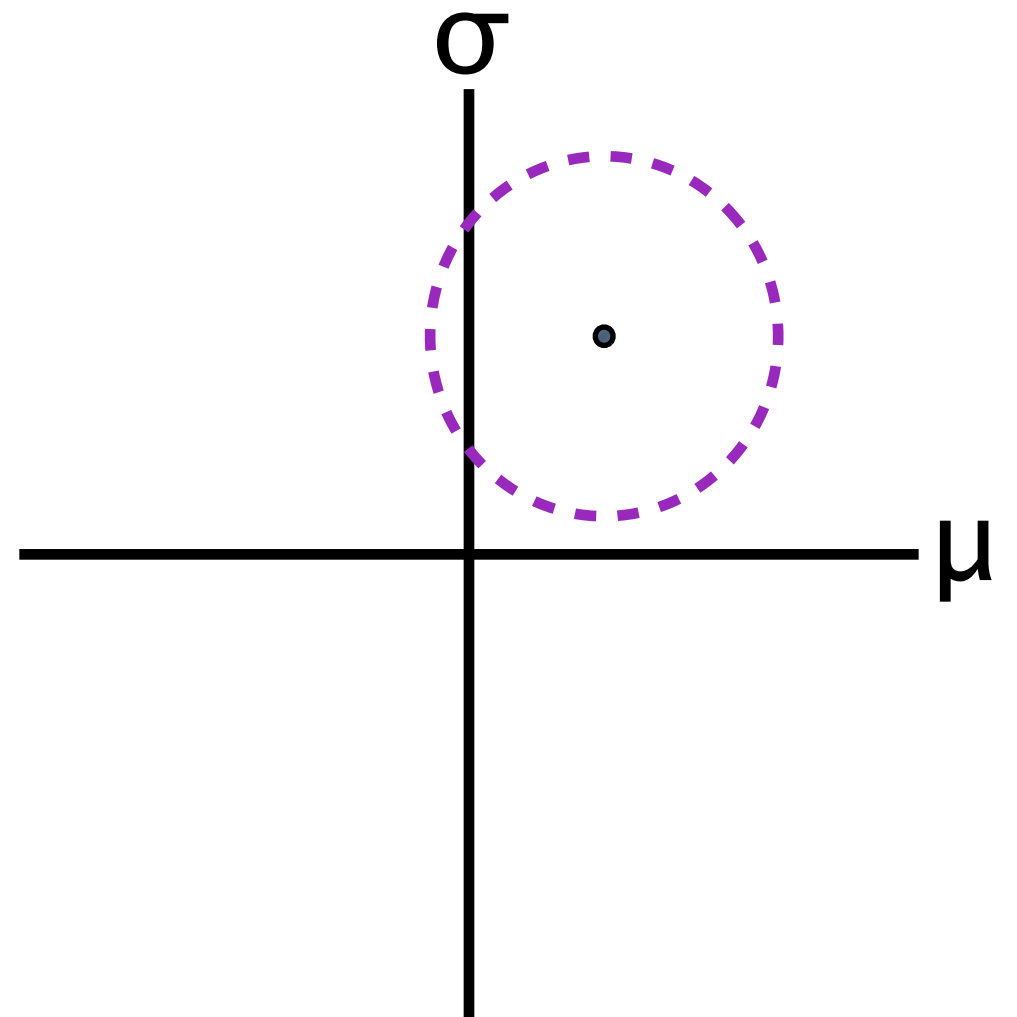


Proximity Variational Inference



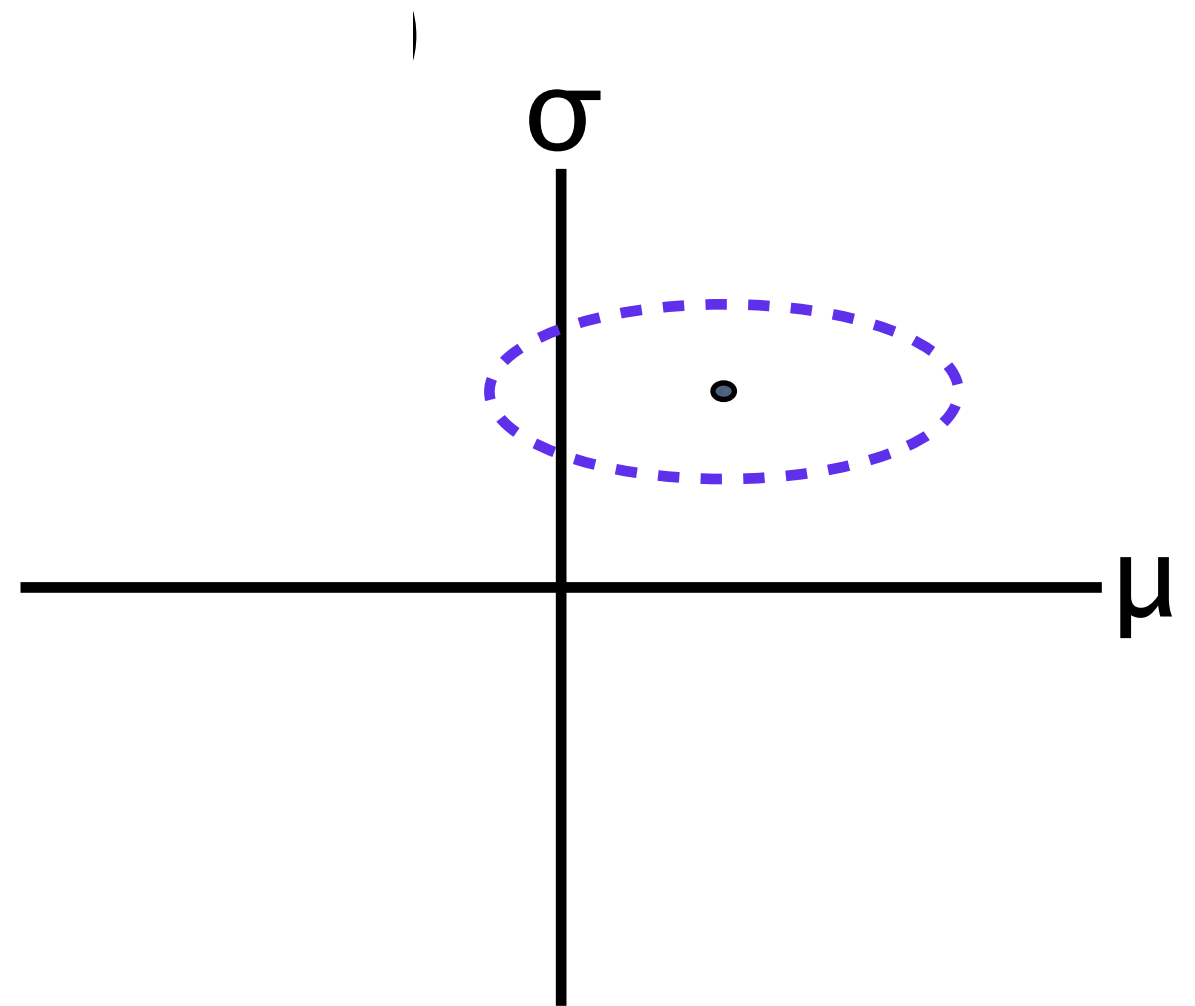
Gradient Ascent: Euclidean Proximity Constraint

$$U(\lambda_{t+1}) = \mathcal{L}(\lambda_t)$$



Proximity Variational Inference

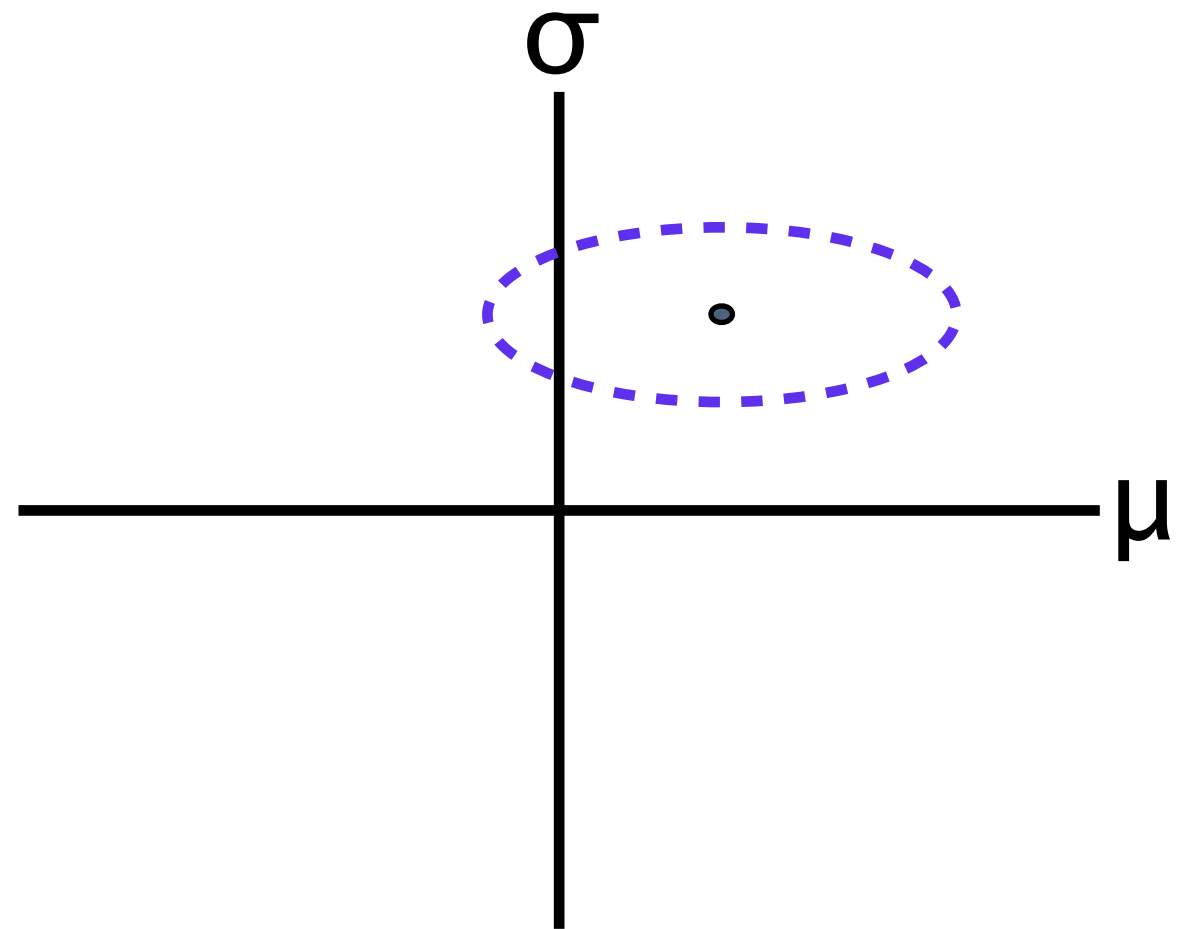
$$U(\boldsymbol{\lambda}_{t+1}) = \mathcal{L}(\boldsymbol{\lambda}_t) + \nabla \mathcal{L}(\boldsymbol{\lambda}_t)^\top (\boldsymbol{\lambda}_{t+1} - \boldsymbol{\lambda}_t) - \frac{1}{2\rho} (\boldsymbol{\lambda}_{t+1} - \boldsymbol{\lambda}_t)^\top (\boldsymbol{\lambda}_{t+1} - \boldsymbol{\lambda}_t)$$

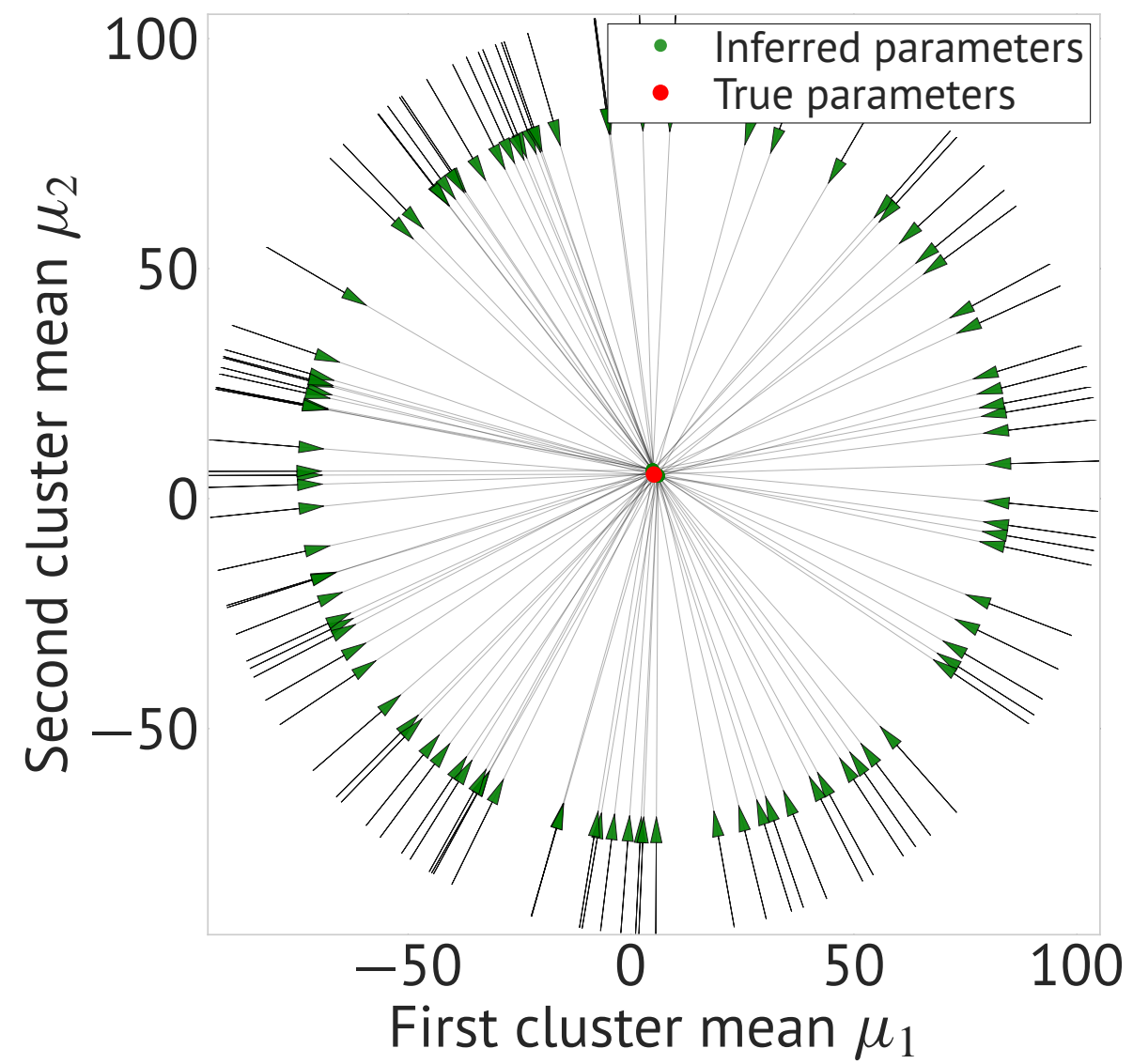
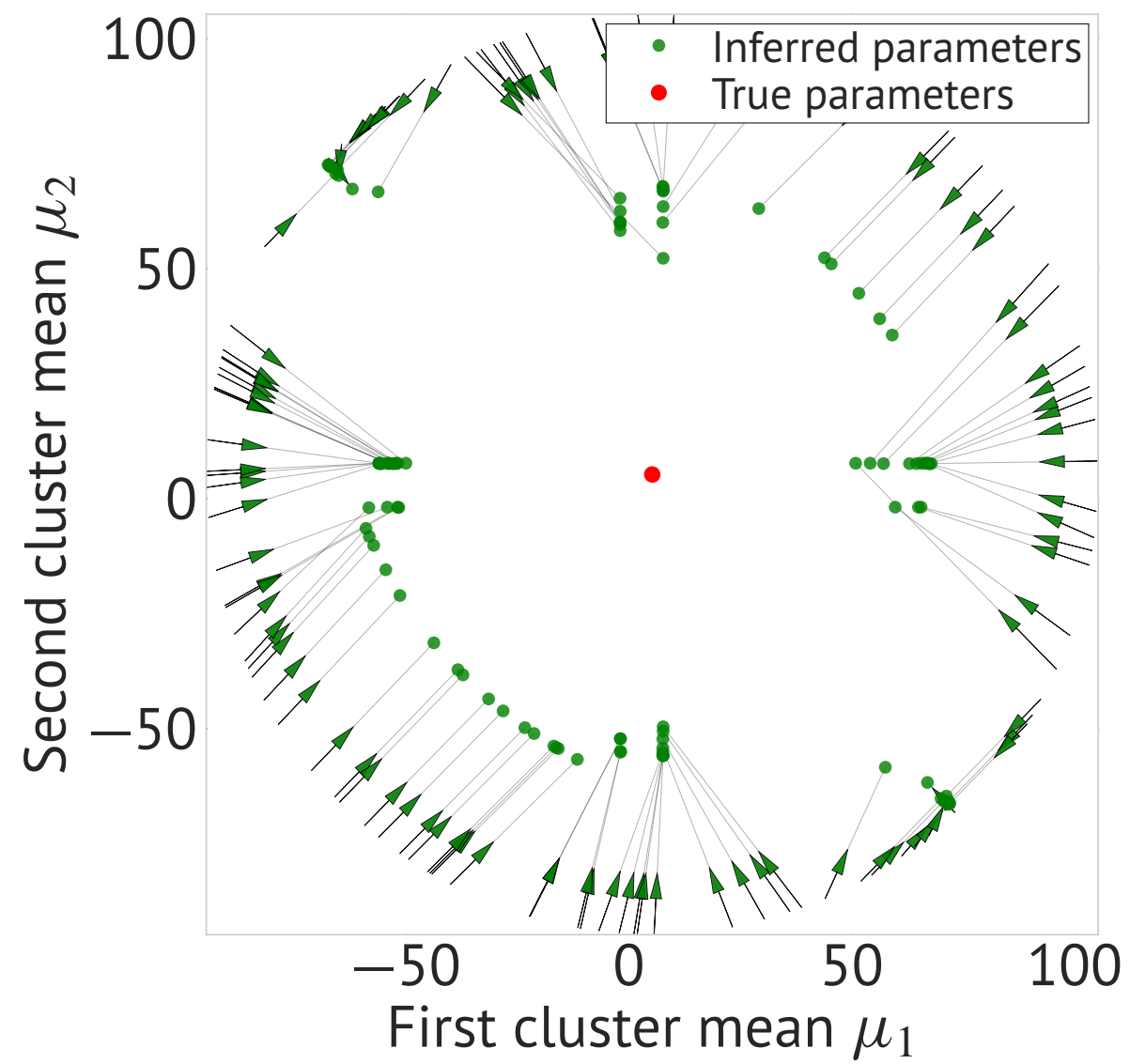


Proximity Variational Inference

Examples of proximity statistics $f(\lambda)$:

- Mean/Variance
- Entropy
- KL divergence
- Many more!



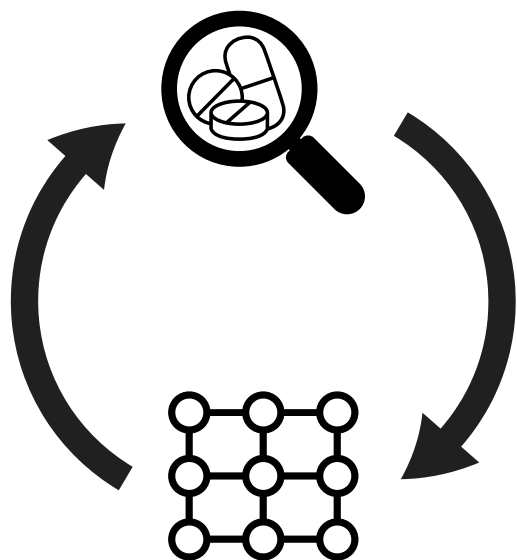


Sigmoid Belief Network

Inference Method	ELBO	Likelihood
Variational Inference	-121.4	-113.7
Deterministic Annealing	-116.8	-108.8
PVI, Entropy Constraint	-113.3	-106.7
PVI, Mean/Variance Constraint	-114.9	-107.4

Ising Model

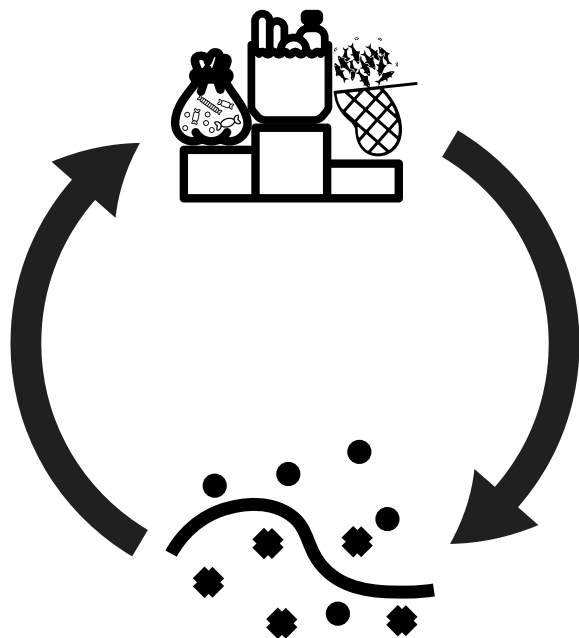
Inference Method	Free Energy
Variational Inference	-2.144
PVI, Entropy Constraint	-2.158



RankFromSets:

Meal Recommendation

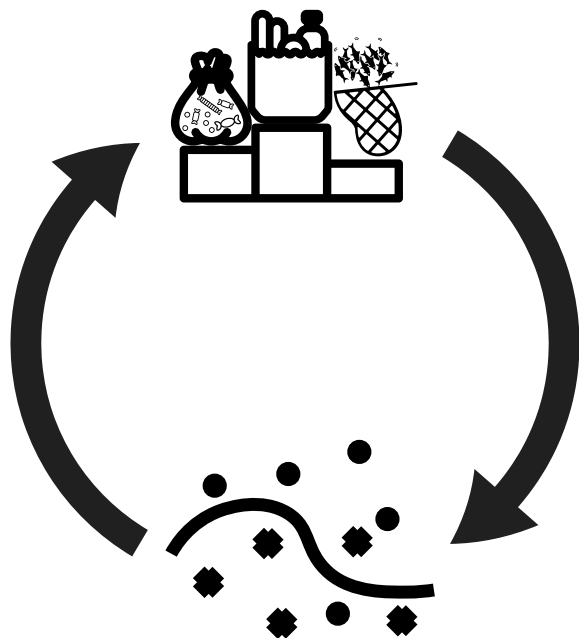
Model	Sampled Recall (%)
RFS	58
RFS, Entropy Constraint	62

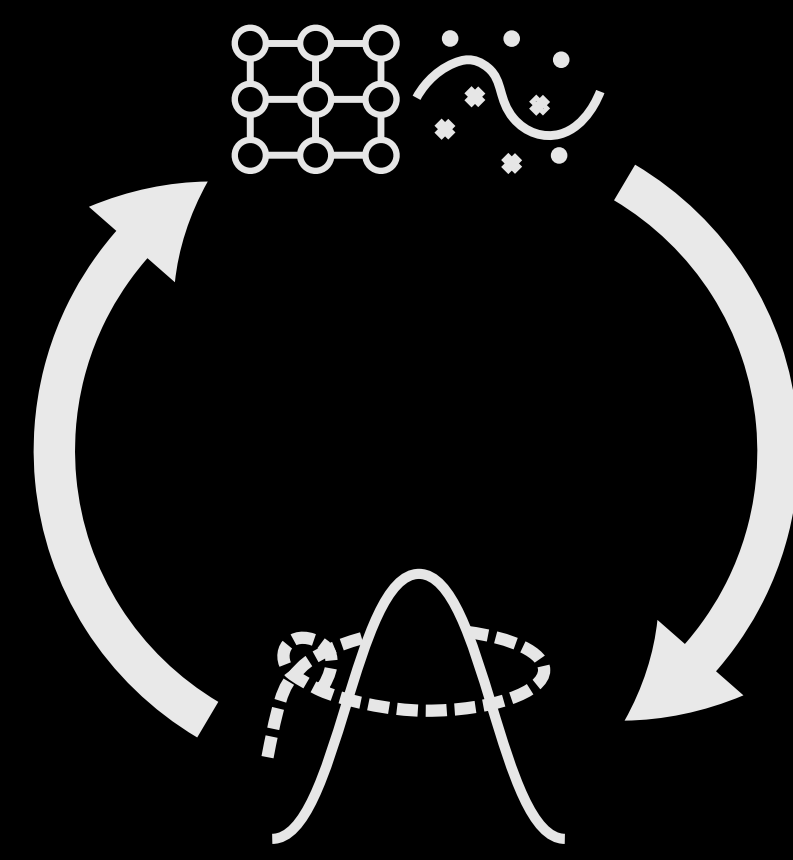
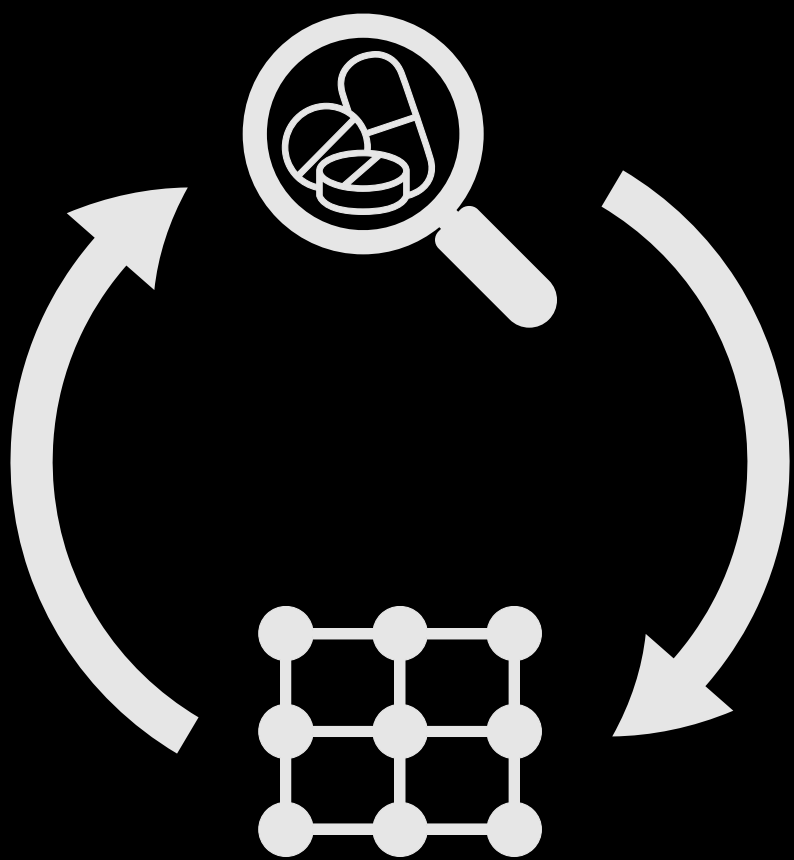


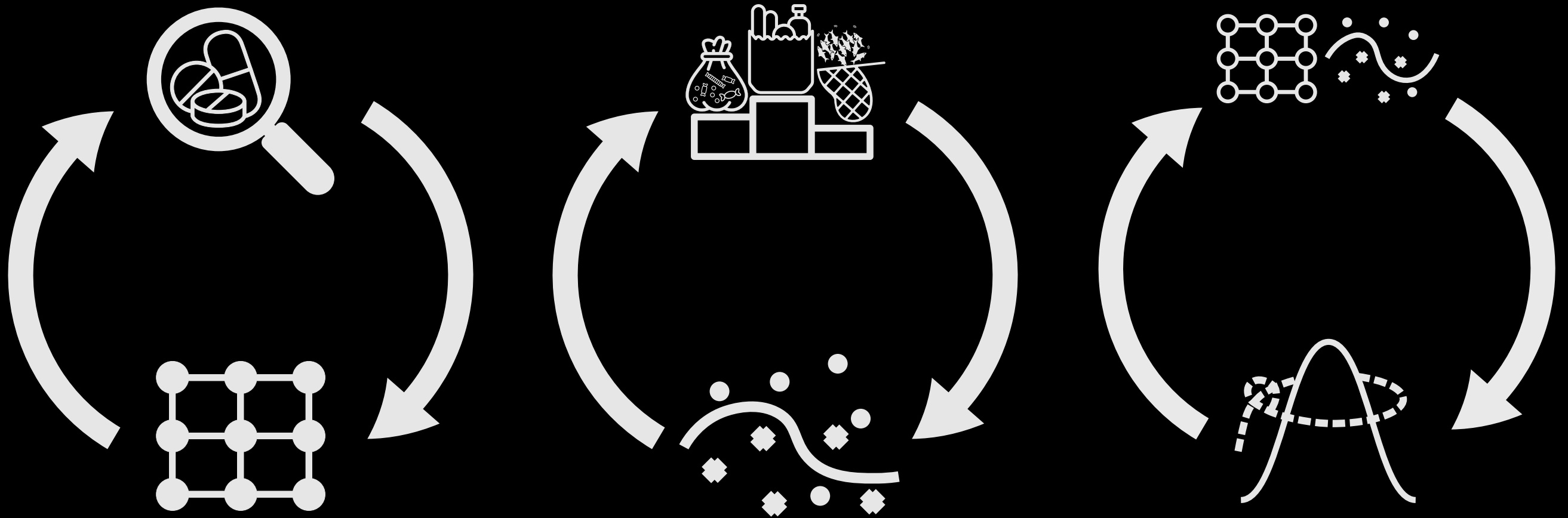
RankFromSets:

ArXiv Recommendation

Model	Recall @ 10 (%)	Recall @ 100 (%)
RFS	0.32	2.54
RFS, Entropy Constraint	0.44	2.54







- J. Altosaar, R. Ranganath, and K. Cranmer. Hierarchical variational models for statistical physics. *MACHINE LEARNING AND THE PHYSICAL SCIENCES WORKSHOP, NEURAL INFORMATION PROCESSING SYSTEMS*. 2019
- J. Altosaar, R. Ranganath, W. Tansey. RankFromSets: Scalable Set Recommendation with Optimal Recall. *AMERICAN STATISTICAL ASSOCIATION SDSS 2020*.
- Altosaar, Jaan, Rajesh Ranganath, and David Blei (2018). "Proximity Variational Inference". *ARTIFICIAL INTELLIGENCE AND STATISTICS*.
- *DESIGN CREDITS:* Sergey Demushkin, Creaticca Creative Agency, Made by Made, Gan Khoun Lay, myiconfinder, Susannanova, Walmart, Raisin Bran, PNGFuel, Alibaba Group, Keith Haring