

Matching

Lecture 7: CV models

Nikita Bozhedomov

Moscow Institute of Physics and Technology

Autumn 2023

Plan

- Types of pictures
- CV in the Matching pipeline
- Past and current research

Types of pictures

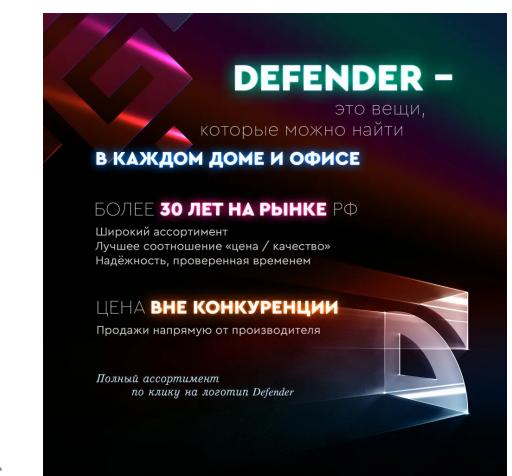
Images



Trash Predictor

Clear

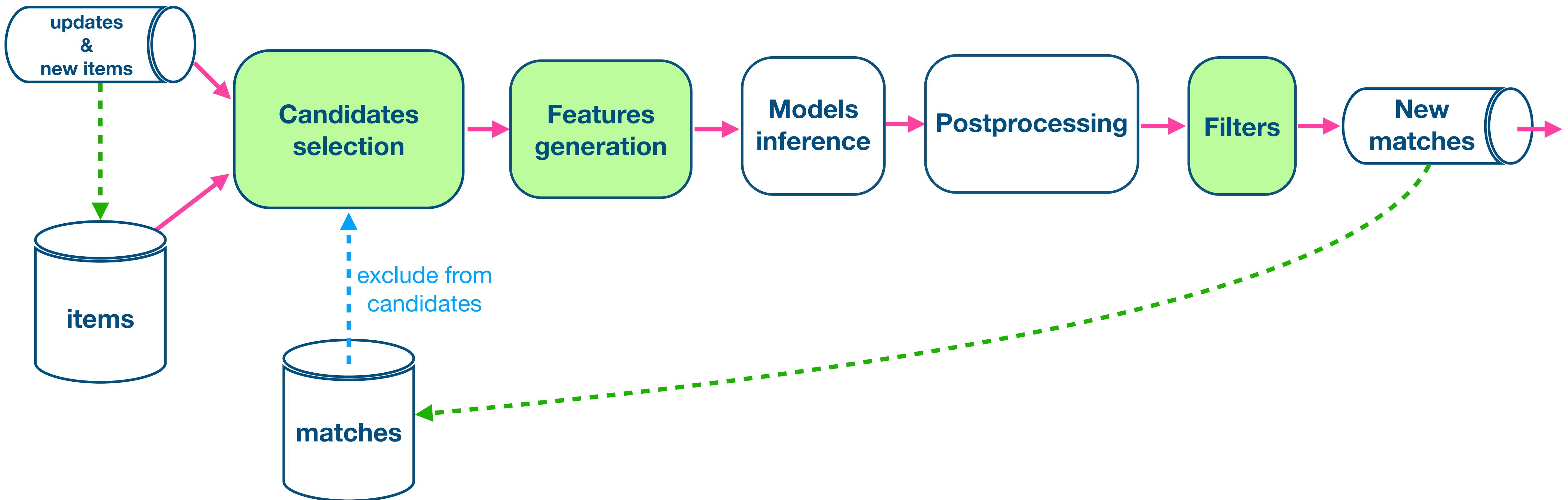
Trash



16,5 см
9,5 см

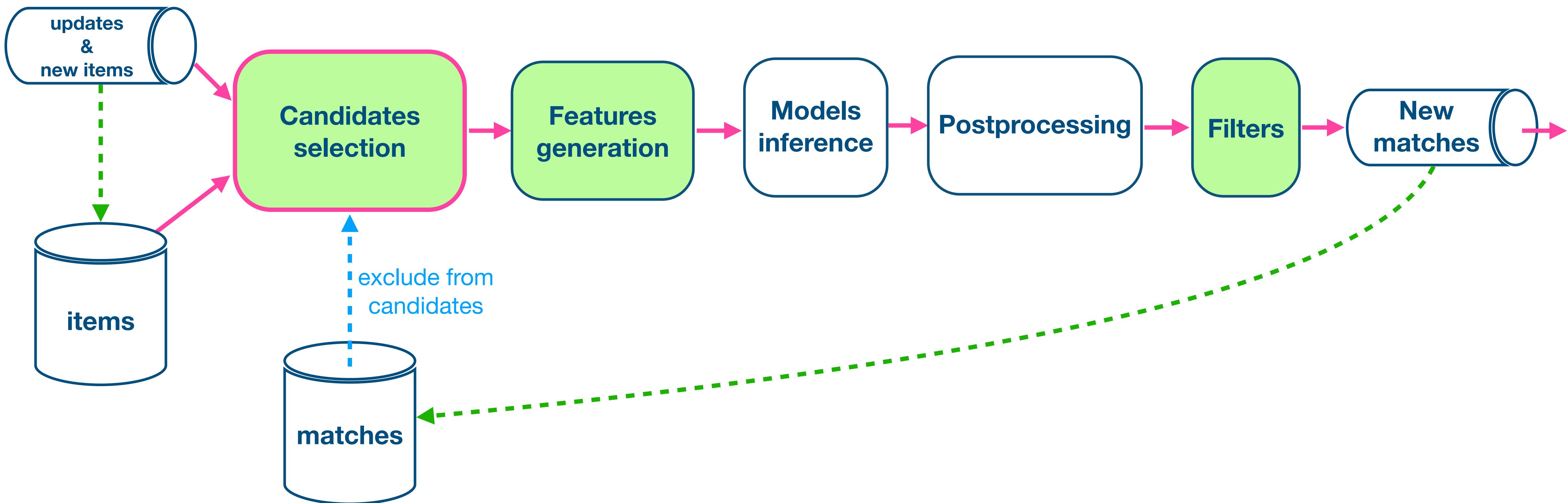
Matching Pipeline

High level design

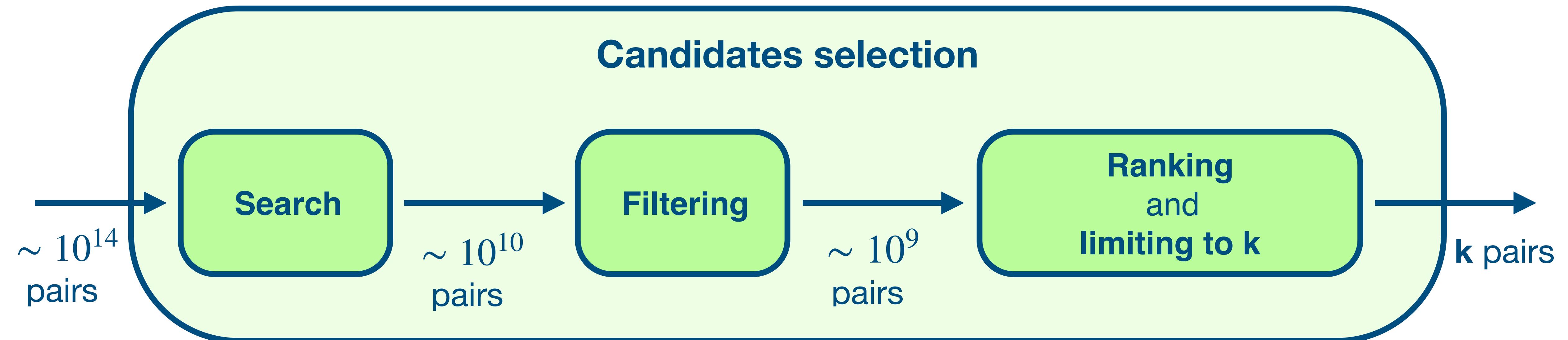


Matching Pipeline

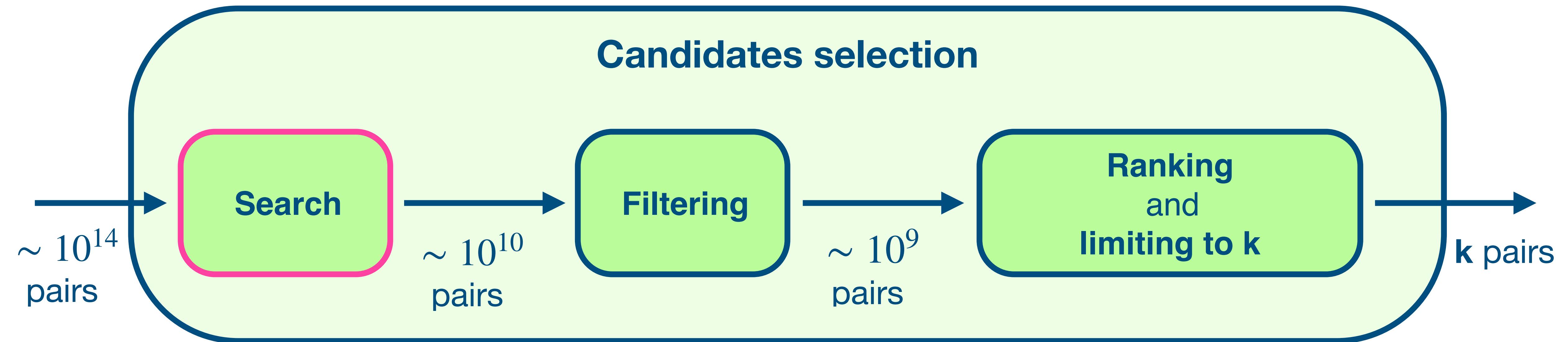
High level design



Candidates selection stages

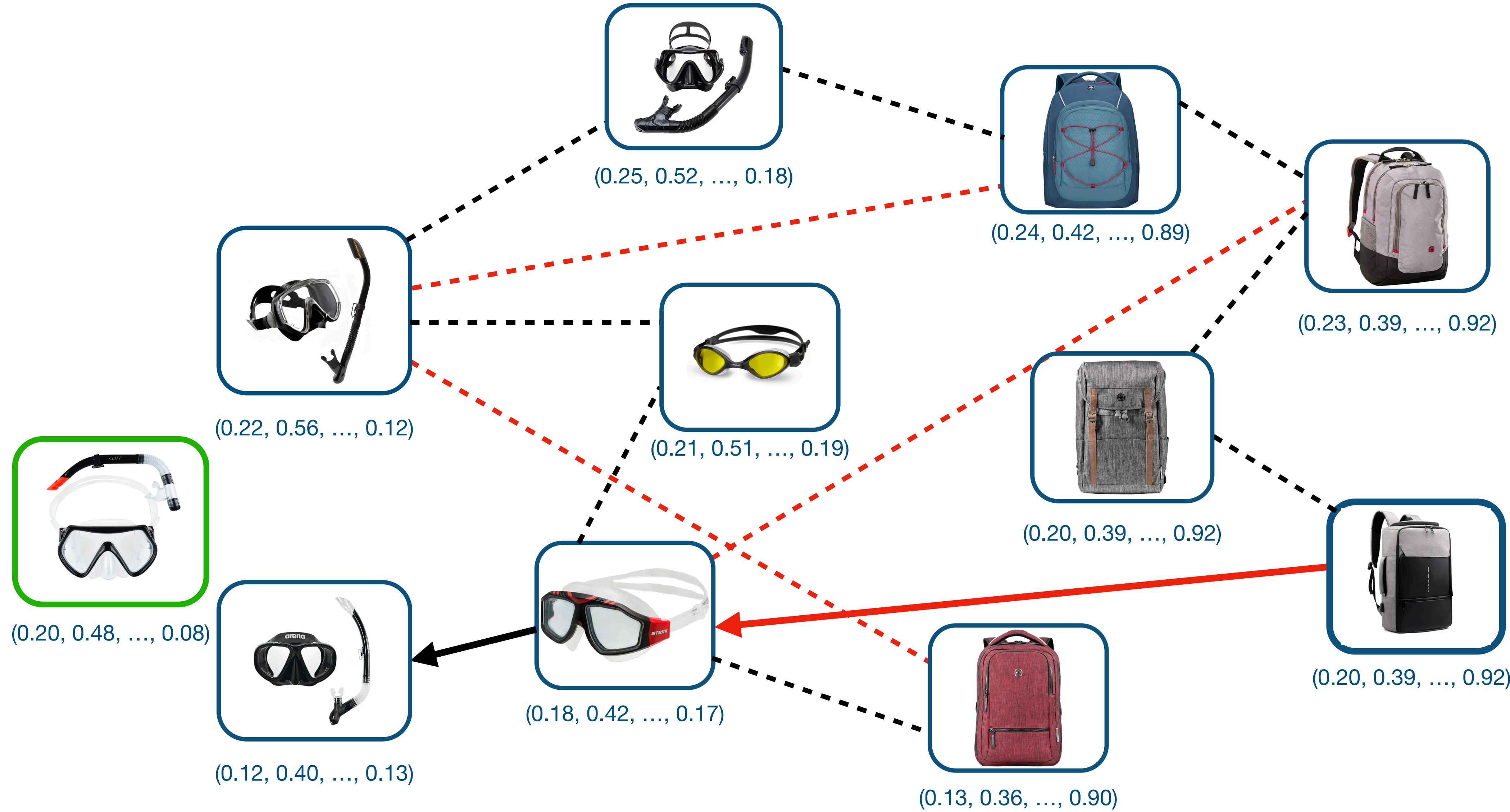


Candidates selection stages

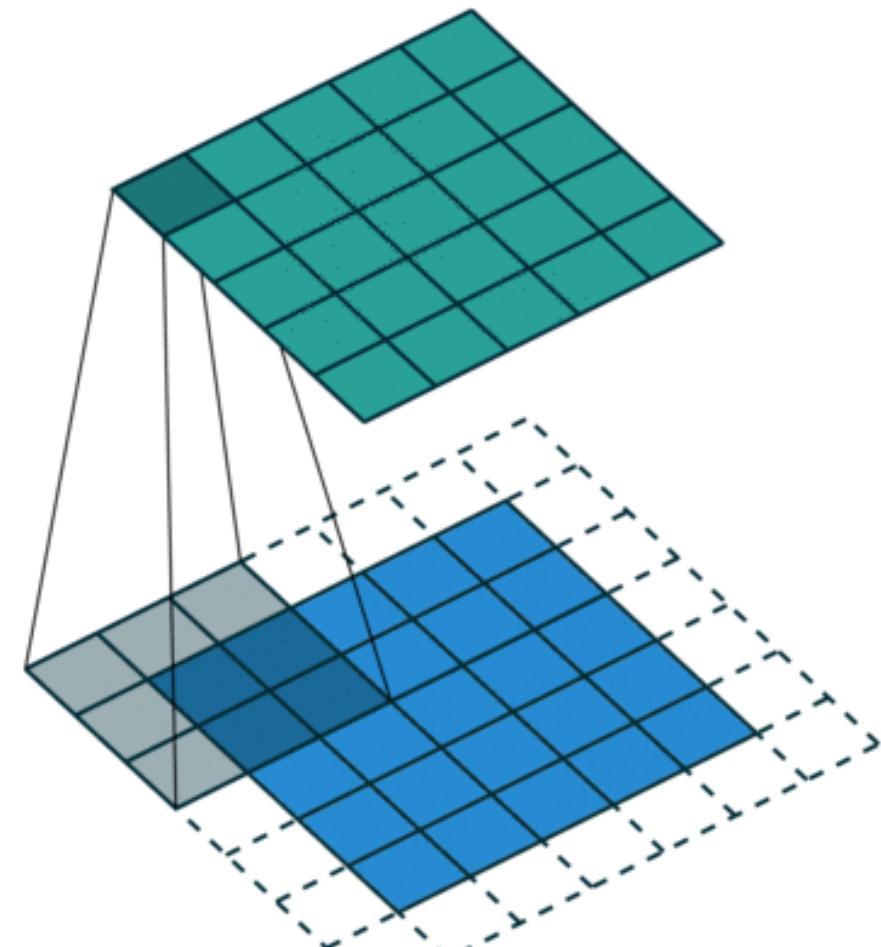


Search

HNSW (Hierarchical Navigable Small World)



Neural Networks



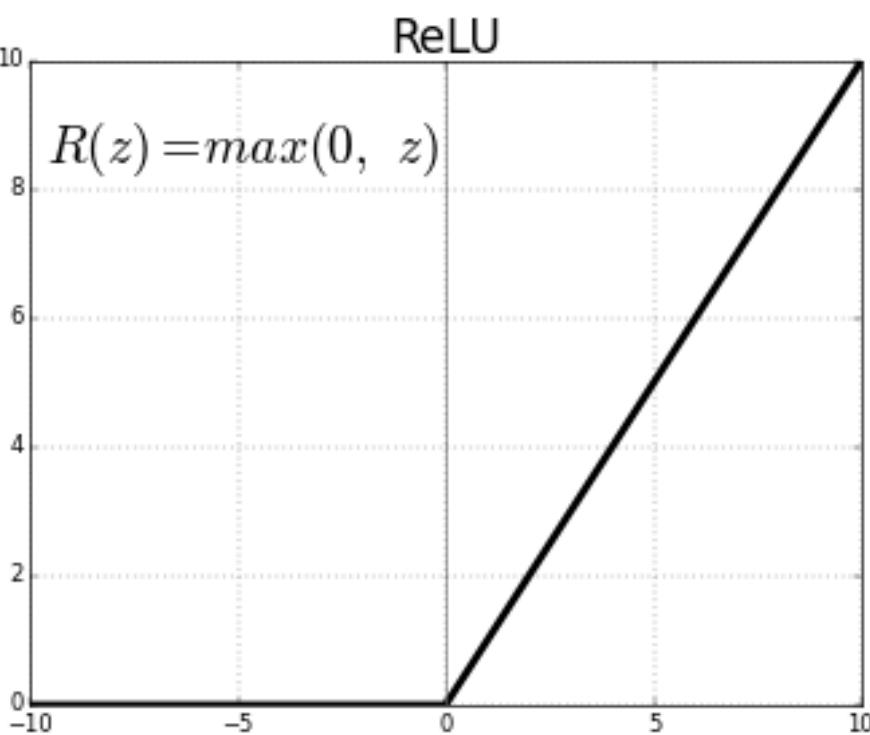
Convolution

$$\begin{aligned}\mu_{\mathcal{B}} &\leftarrow \frac{1}{m} \sum_{i=1}^m x_i \\ \sigma_{\mathcal{B}}^2 &\leftarrow \frac{1}{m} \sum_{i=1}^m (x_i - \mu_{\mathcal{B}})^2 \\ \hat{x}_i &\leftarrow \frac{x_i - \mu_{\mathcal{B}}}{\sqrt{\sigma_{\mathcal{B}}^2 + \epsilon}} \\ y_i &\leftarrow \gamma \hat{x}_i + \beta \equiv \text{BN}_{\gamma, \beta}(x_i)\end{aligned}$$

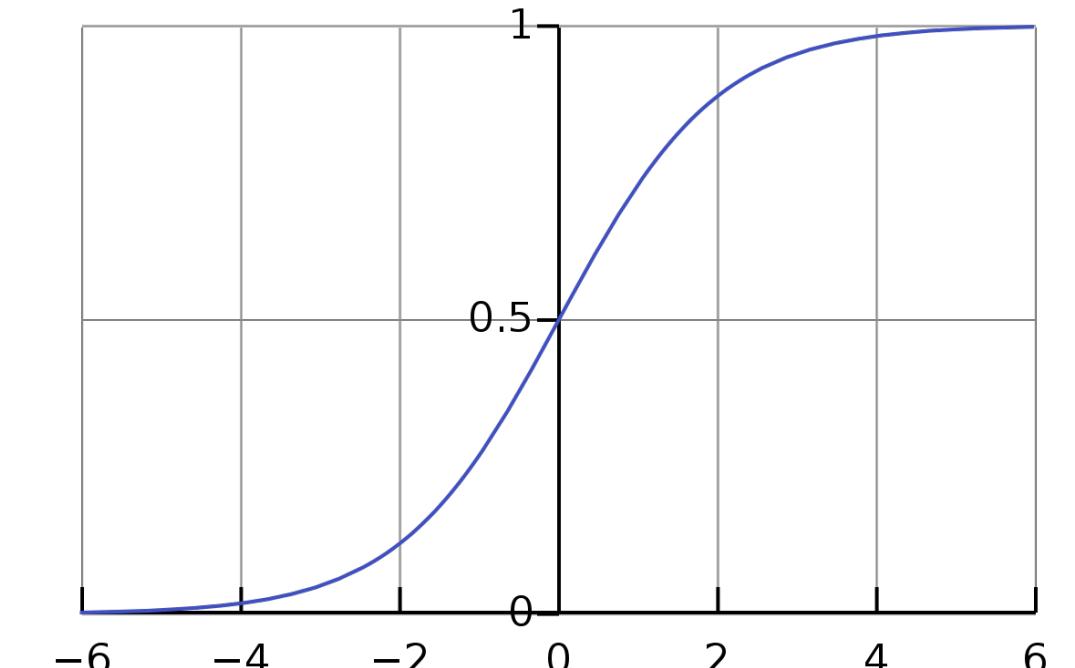
Batch Normalization

[S. Ioffe , C. Szegedy, 2015]

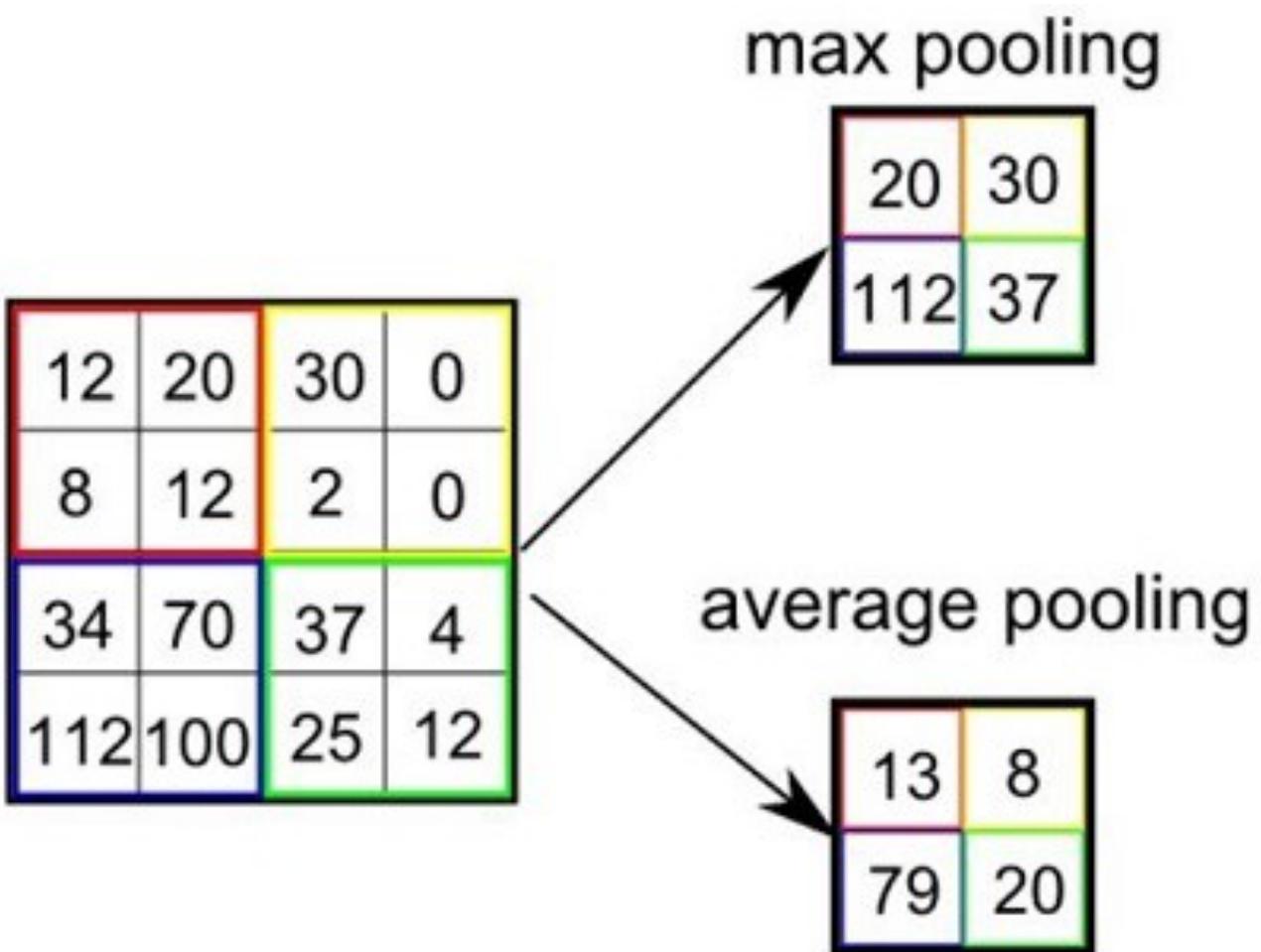
Base layers



Nonlinear activation: ReLU



Nonlinear activation: Sigmoid



Pooling

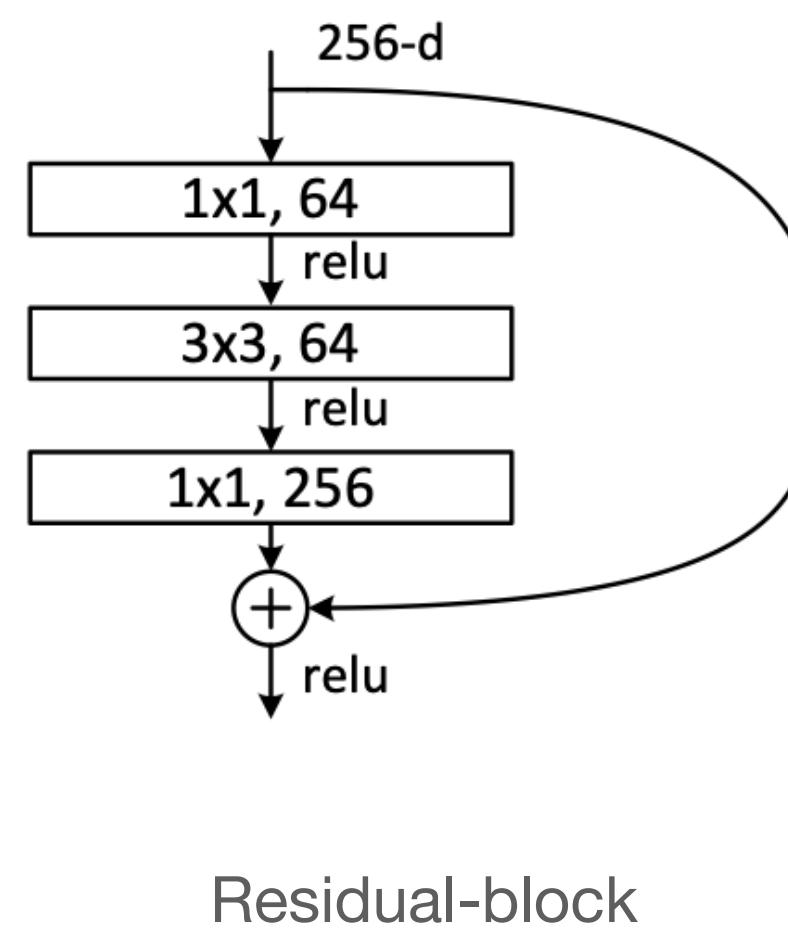
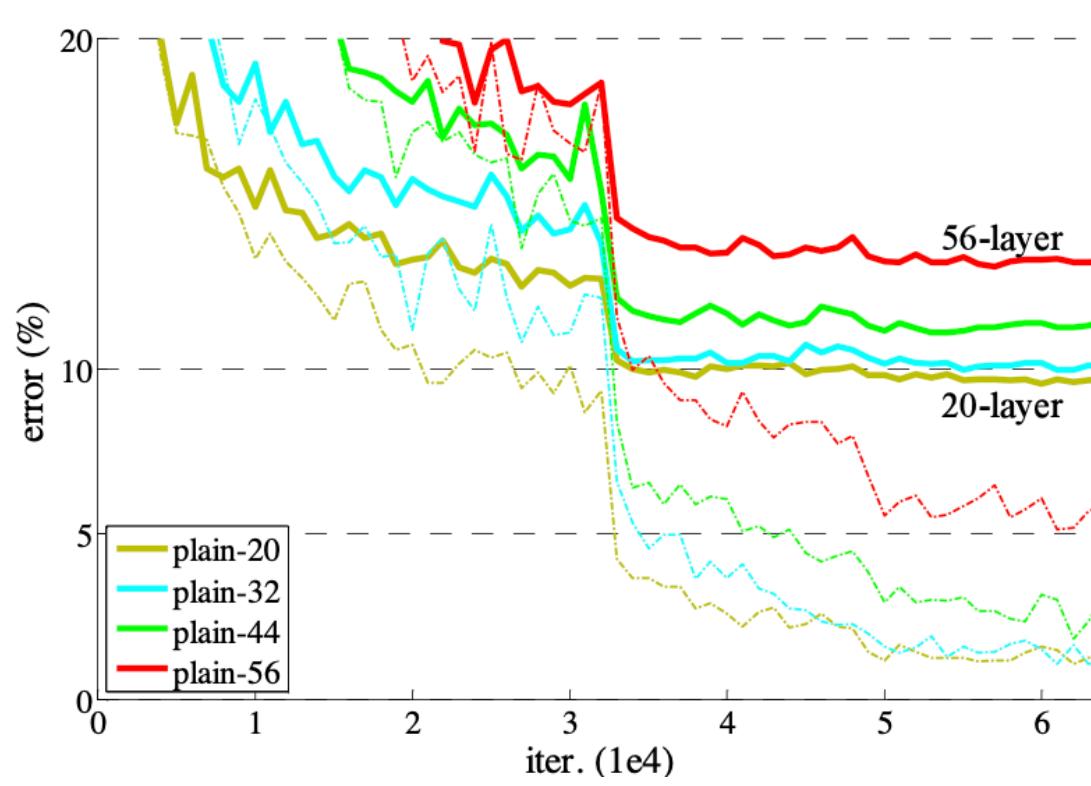
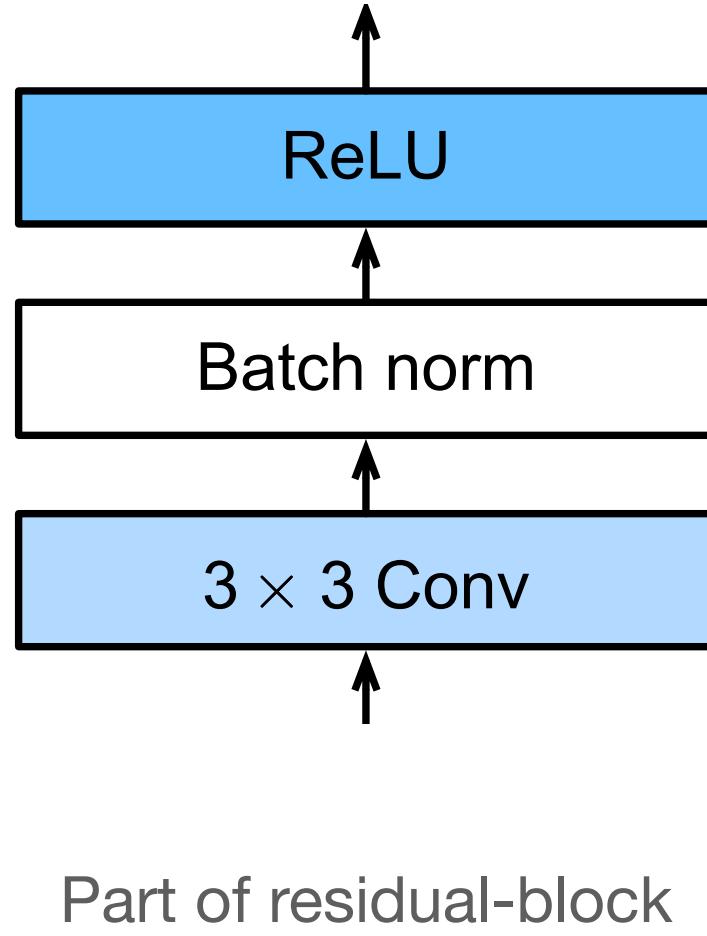
$$\sigma : \mathbb{R}^K \mapsto (0, 1)^K$$

$$\sigma(\mathbf{z})_i = \frac{e^{z_i}}{\sum_{j=1}^K e^{z_j}}$$

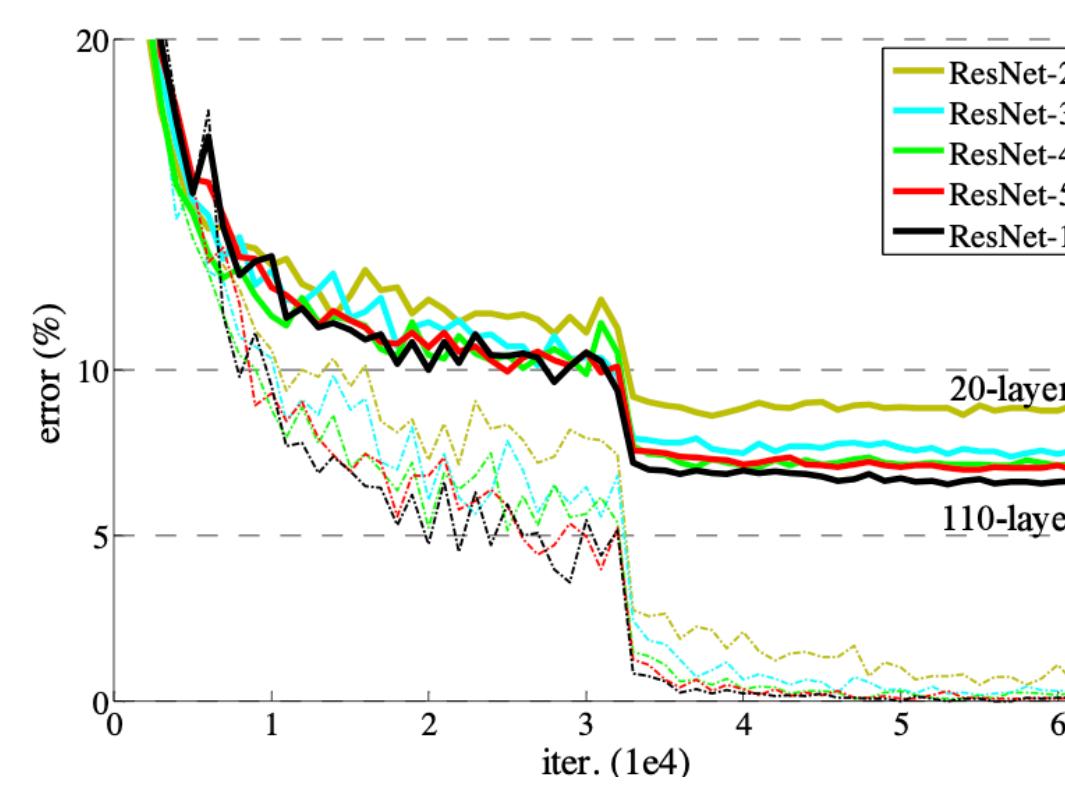
Softmax

Neural Networks

ResNet. Architecture



layer name	output size	18-layer	34-layer	50-layer	101-layer	152-layer
conv1	112×112			7×7, 64, stride 2		
conv2_x	56×56	$\left[\begin{matrix} 3\times 3, 64 \\ 3\times 3, 64 \end{matrix} \right] \times 2$	$\left[\begin{matrix} 3\times 3, 64 \\ 3\times 3, 64 \end{matrix} \right] \times 3$	$\left[\begin{matrix} 1\times 1, 64 \\ 3\times 3, 64 \\ 1\times 1, 256 \end{matrix} \right] \times 3$	$\left[\begin{matrix} 1\times 1, 64 \\ 3\times 3, 64 \\ 1\times 1, 256 \end{matrix} \right] \times 3$	$\left[\begin{matrix} 1\times 1, 64 \\ 3\times 3, 64 \\ 1\times 1, 256 \end{matrix} \right] \times 3$
conv3_x	28×28	$\left[\begin{matrix} 3\times 3, 128 \\ 3\times 3, 128 \end{matrix} \right] \times 2$	$\left[\begin{matrix} 3\times 3, 128 \\ 3\times 3, 128 \end{matrix} \right] \times 4$	$\left[\begin{matrix} 1\times 1, 128 \\ 3\times 3, 128 \\ 1\times 1, 512 \end{matrix} \right] \times 4$	$\left[\begin{matrix} 1\times 1, 128 \\ 3\times 3, 128 \\ 1\times 1, 512 \end{matrix} \right] \times 4$	$\left[\begin{matrix} 1\times 1, 128 \\ 3\times 3, 128 \\ 1\times 1, 512 \end{matrix} \right] \times 8$
conv4_x	14×14	$\left[\begin{matrix} 3\times 3, 256 \\ 3\times 3, 256 \end{matrix} \right] \times 2$	$\left[\begin{matrix} 3\times 3, 256 \\ 3\times 3, 256 \end{matrix} \right] \times 6$	$\left[\begin{matrix} 1\times 1, 256 \\ 3\times 3, 256 \\ 1\times 1, 1024 \end{matrix} \right] \times 6$	$\left[\begin{matrix} 1\times 1, 256 \\ 3\times 3, 256 \\ 1\times 1, 1024 \end{matrix} \right] \times 23$	$\left[\begin{matrix} 1\times 1, 256 \\ 3\times 3, 256 \\ 1\times 1, 1024 \end{matrix} \right] \times 36$
conv5_x	7×7	$\left[\begin{matrix} 3\times 3, 512 \\ 3\times 3, 512 \end{matrix} \right] \times 2$	$\left[\begin{matrix} 3\times 3, 512 \\ 3\times 3, 512 \end{matrix} \right] \times 3$	$\left[\begin{matrix} 1\times 1, 512 \\ 3\times 3, 512 \\ 1\times 1, 2048 \end{matrix} \right] \times 3$	$\left[\begin{matrix} 1\times 1, 512 \\ 3\times 3, 512 \\ 1\times 1, 2048 \end{matrix} \right] \times 3$	$\left[\begin{matrix} 1\times 1, 512 \\ 3\times 3, 512 \\ 1\times 1, 2048 \end{matrix} \right] \times 3$
	1×1			average pool, 1000-d fc, softmax		
FLOPs		1.8×10^9	3.6×10^9	3.8×10^9	7.6×10^9	11.3×10^9



Architecture

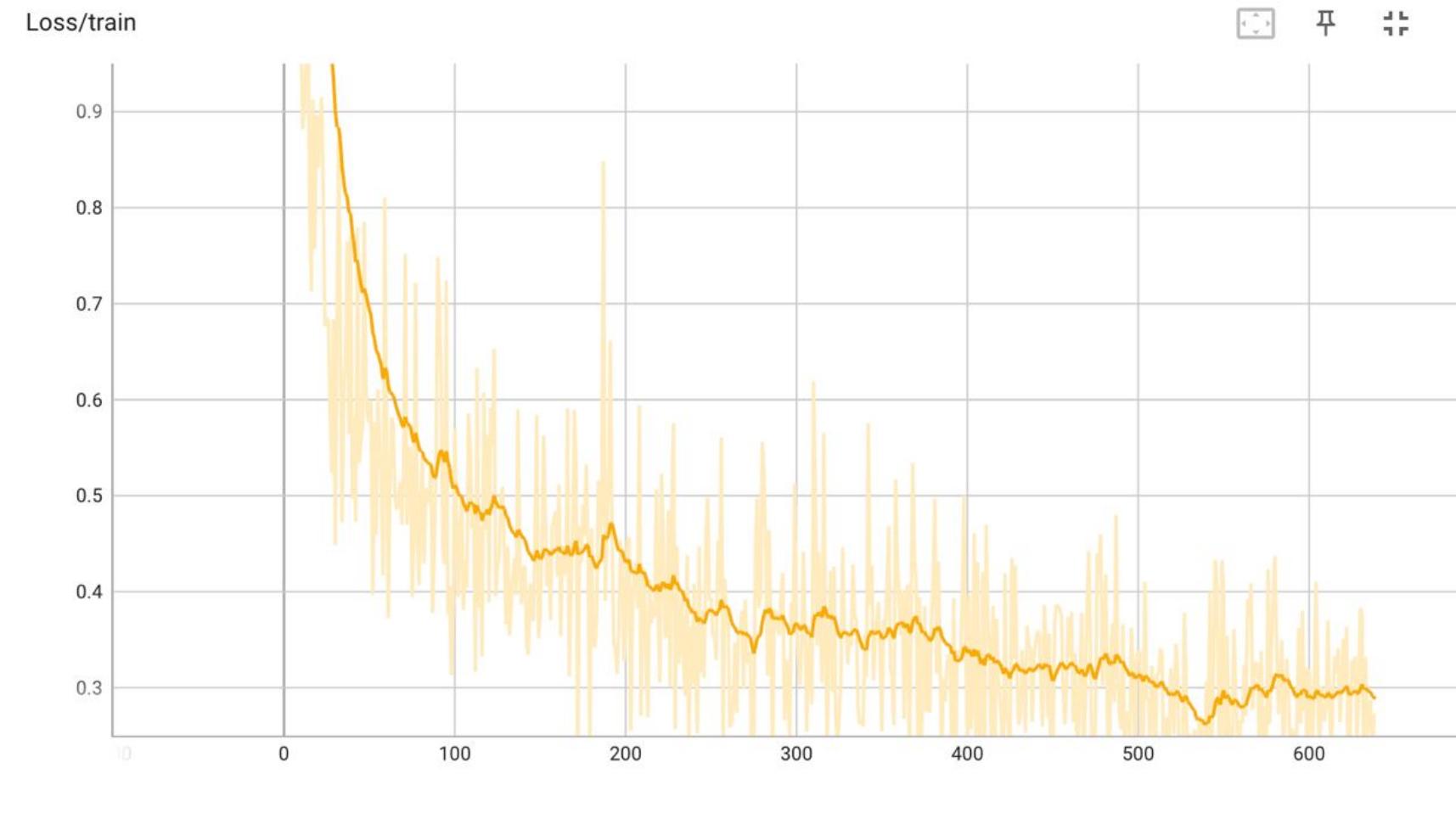
[K. He et al., 2015]

Neural Networks

ResNet. Training

Dataset

	old	new
positives	0.7M	13.3M
negatives		3.9M



NT-Xent loss

[Chen T. et al., 2020]

$$\ell_{i,j} = -\log \frac{\exp(\text{sim}(z_i, z_j)/\tau)}{\sum_{k=1}^{2N} \mathbb{1}_{[k \neq i]} \exp(\text{sim}(z_i, z_k)/\tau)}$$

$$\text{sim}(\mathbf{u}, \mathbf{v}) = \mathbf{u}^T \mathbf{v} / \|\mathbf{u}\| \|\mathbf{v}\|$$

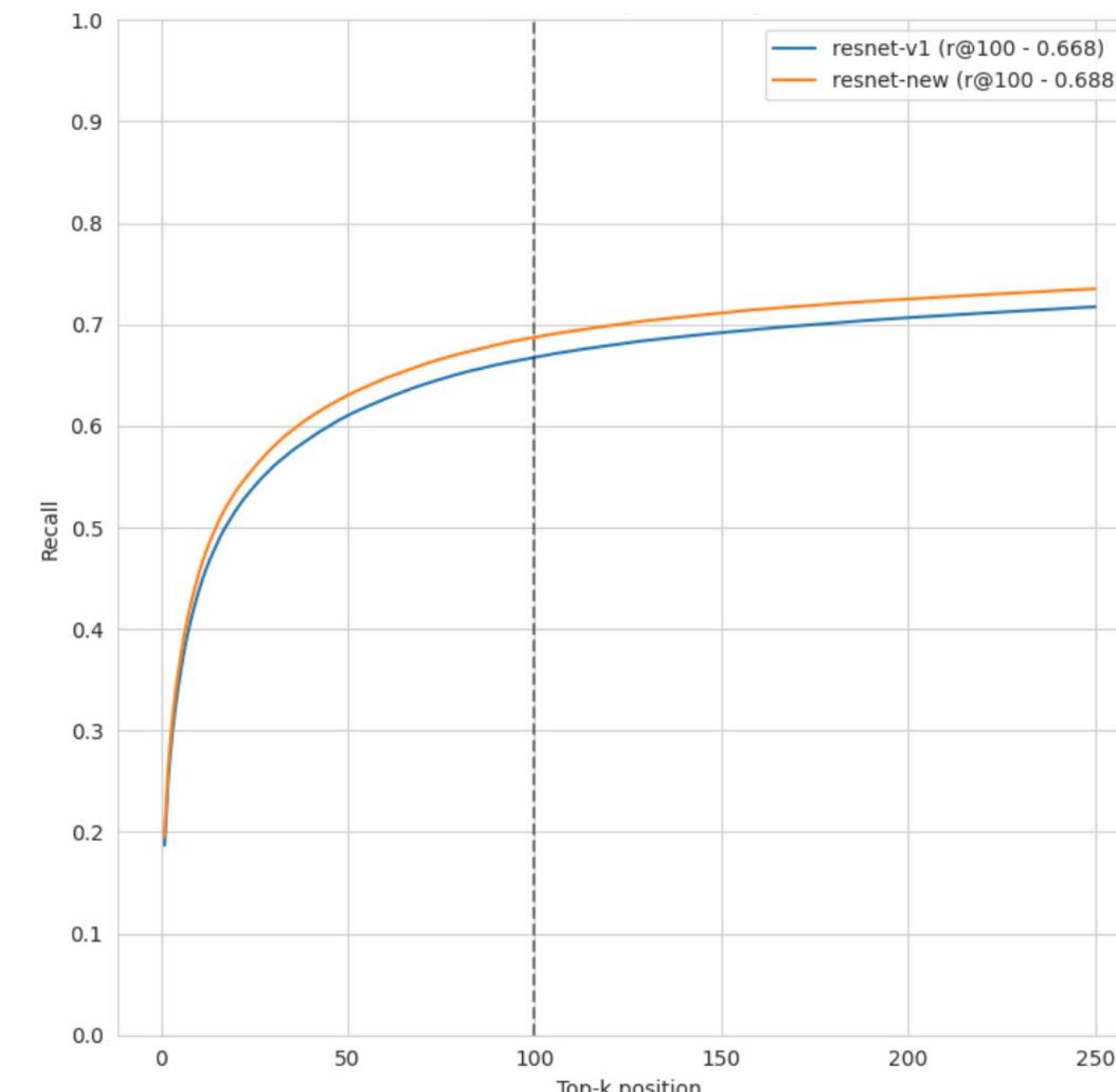
Optimization

Stochastic Gradient Descent

$$w := w - \eta \nabla Q(w) = w - \frac{\eta}{n} \sum_{i=1}^n \nabla Q_i(w)$$

$lr = \eta = 0.01$, $batch_size = n = 128$

$n_process = 4$, $ratio = 1 \frac{GPU}{process}$



Neural Networks

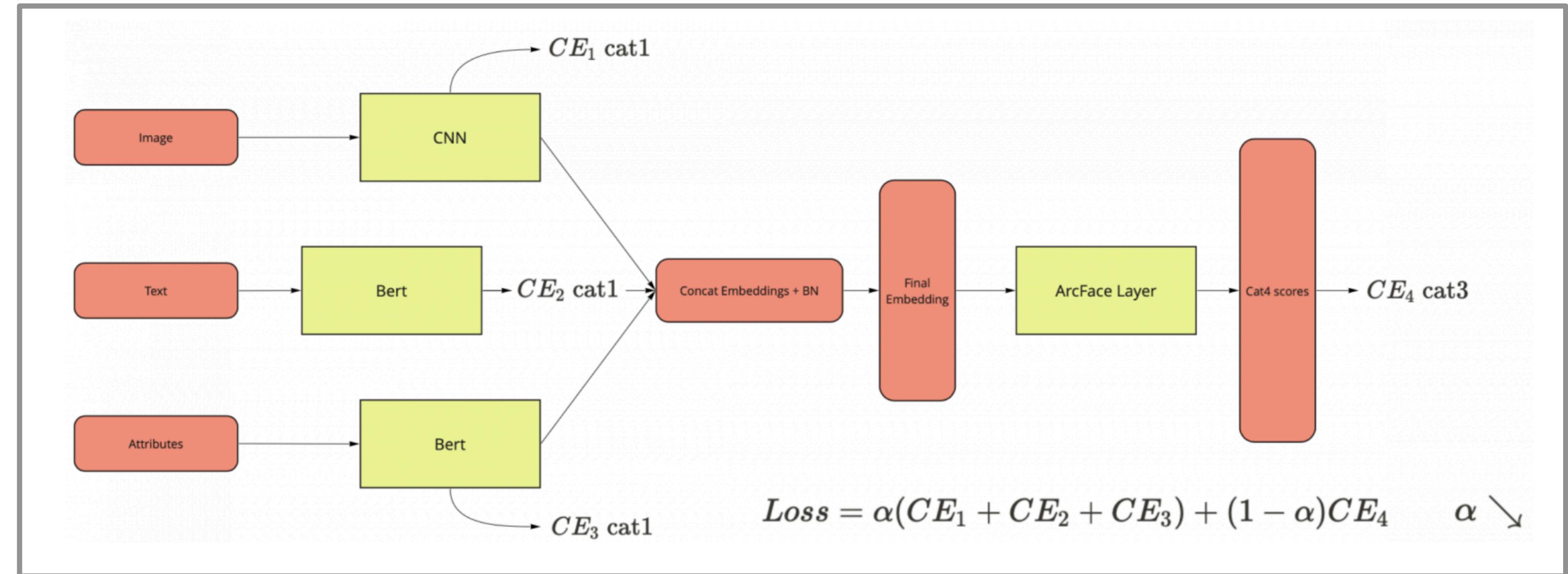
Previous research

$BN \rightarrow AdaptiveGradientClipping$

$$\|W^l\|_F = \sqrt{\sum_{i,j} (W_{i,j}^l)^2} - \text{Frobenius norm}$$

$$G_i^\ell \rightarrow \begin{cases} \lambda \frac{\|W_i^\ell\|_F^*}{\|G_i^\ell\|_F} G_i^\ell & \text{if } \frac{\|G_i^\ell\|_F}{\|W_i^\ell\|_F^*} > \lambda, \\ G_i^\ell & \text{otherwise.} \end{cases}$$

$$\|W_i\|_F^* = \max(\|W_i\|_F, \epsilon)$$



Normalized-Free ResNets

[A. Brock et al., 2021]

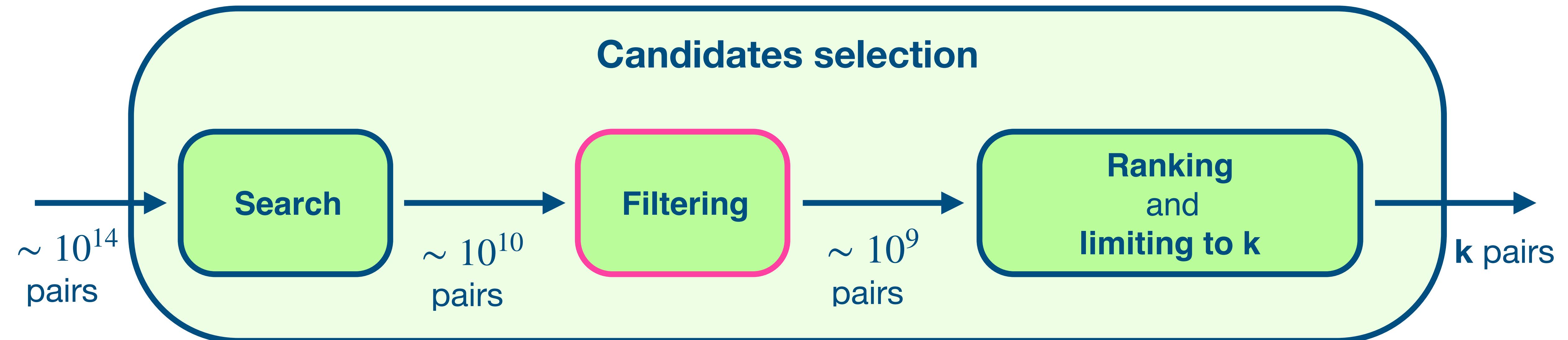
$$L_3 = -\log \frac{e^{s \cos(\theta_{y_i} + m)}}{e^{s \cos(\theta_{y_i} + m)} + \sum_{j=1, j \neq y_i}^N e^{s \cos \theta_j}}$$

ArcFace loss
[J. Deng et al., 2015]

Prod2Vec

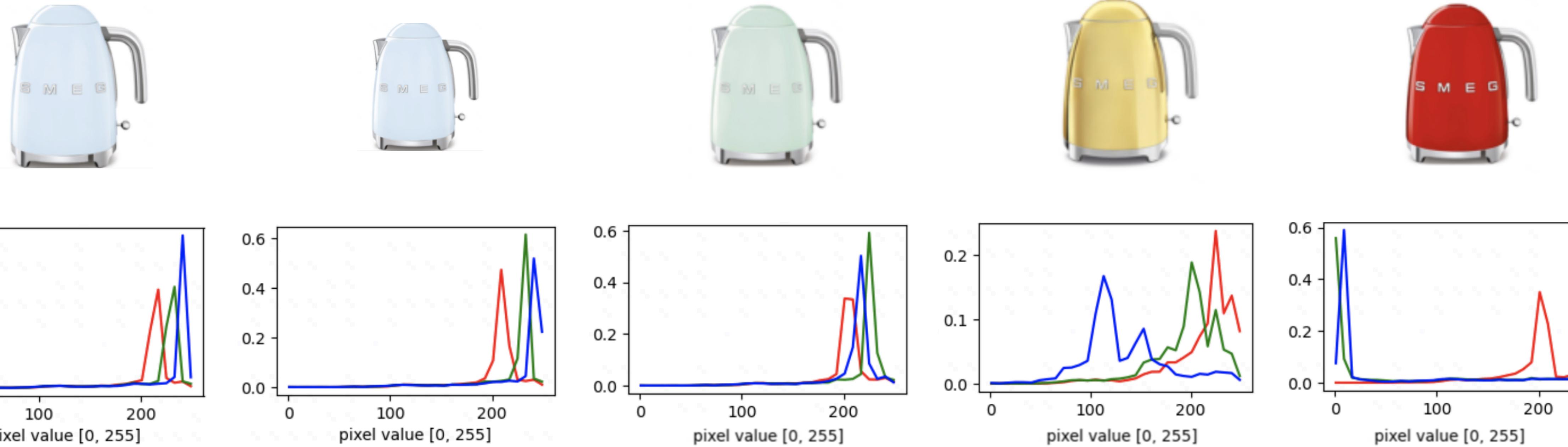
[Ozon Matching, 2021]

Candidates selection stages



Filtering

Filter by color

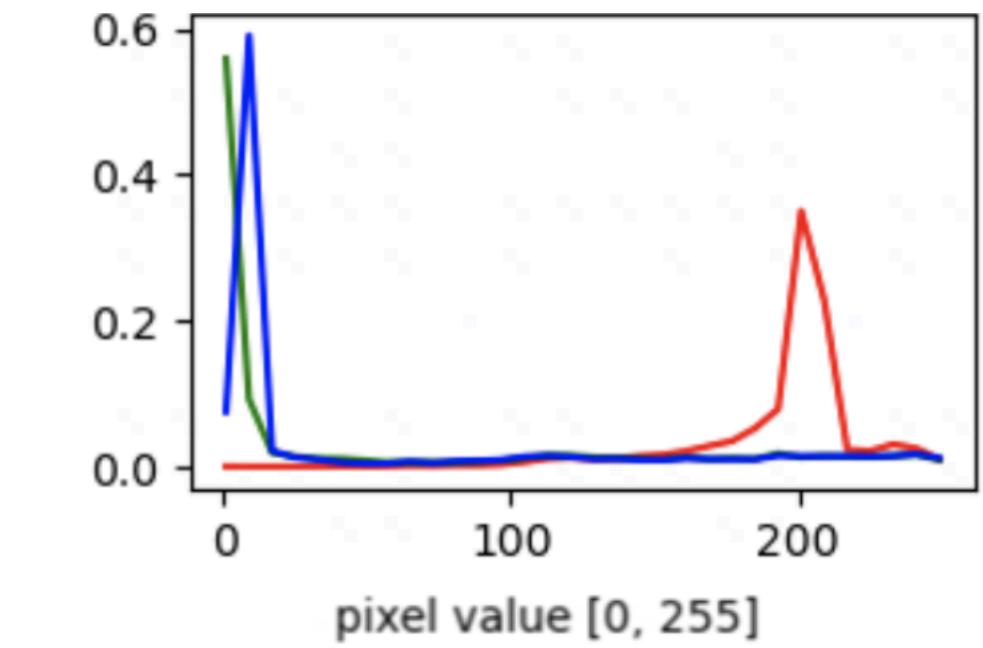
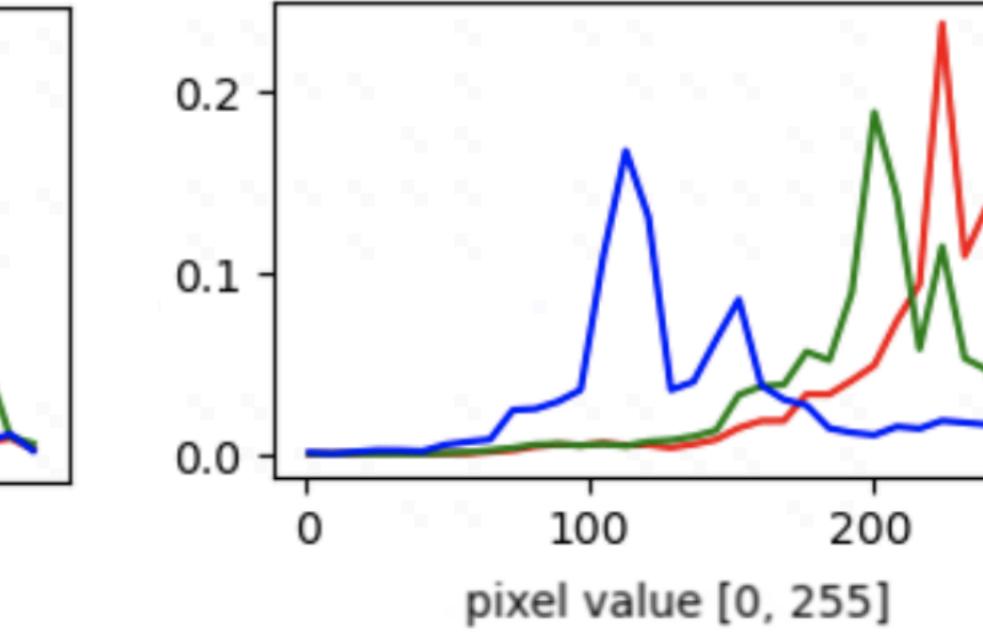
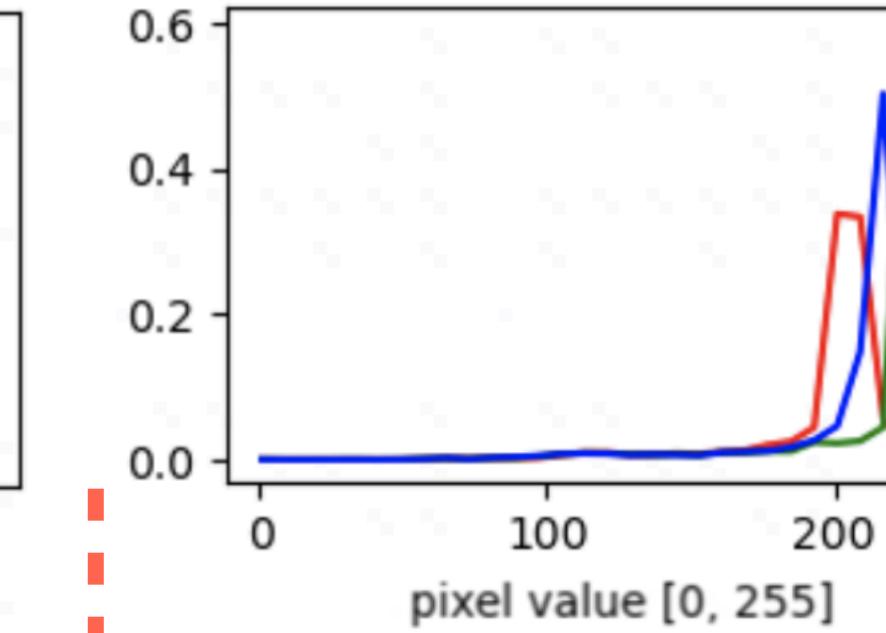
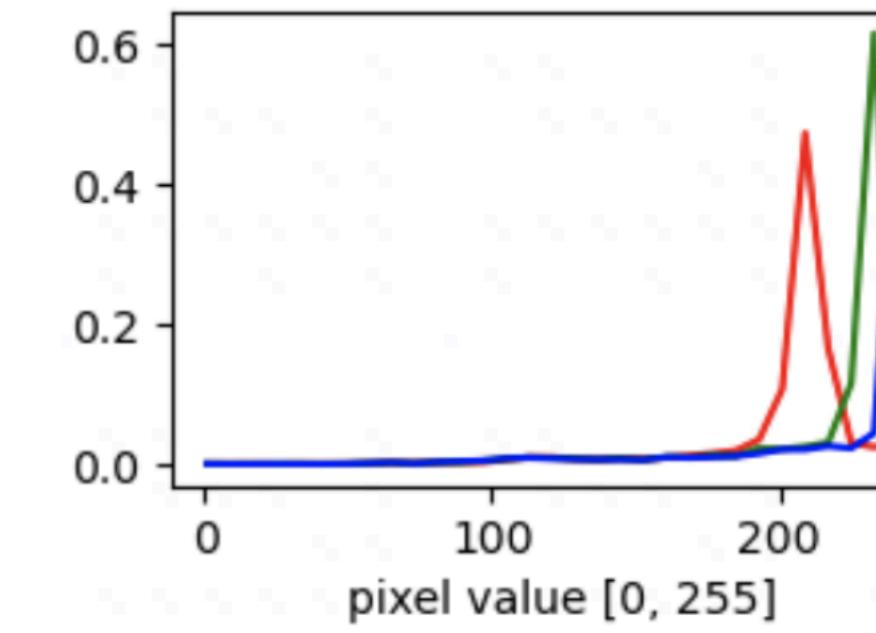
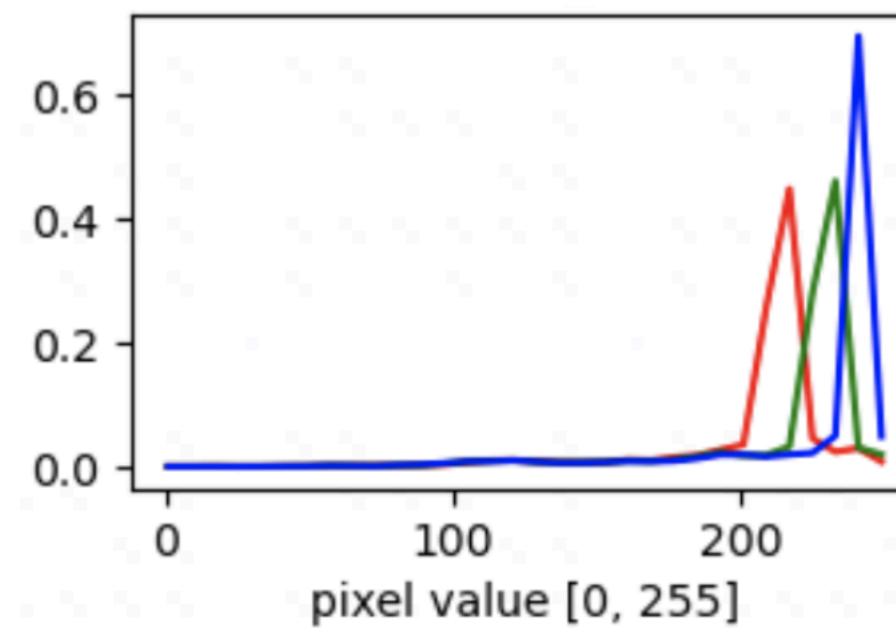


$$\text{Jeffreys divergence} = D_{\text{KL}}(P \parallel Q) + D_{\text{KL}}(Q \parallel P)$$

$$D_{\text{KL}}(P \parallel Q) = \sum_{x \in \mathcal{X}} P(x) \log\left(\frac{P(x)}{Q(x)}\right)$$

Filtering

Filter by color



? Possible match

X Not match

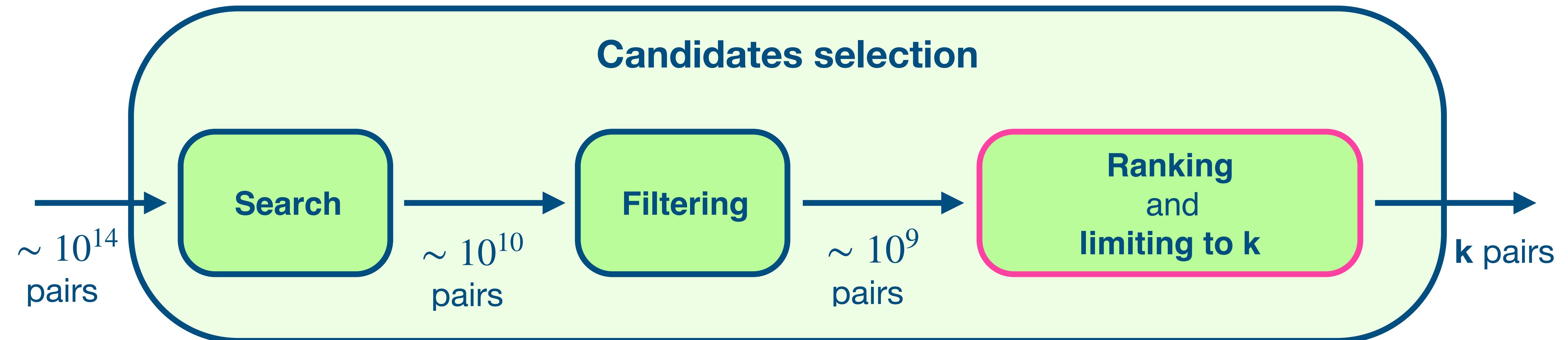
Divergence = 0.34

Divergence = 2.07

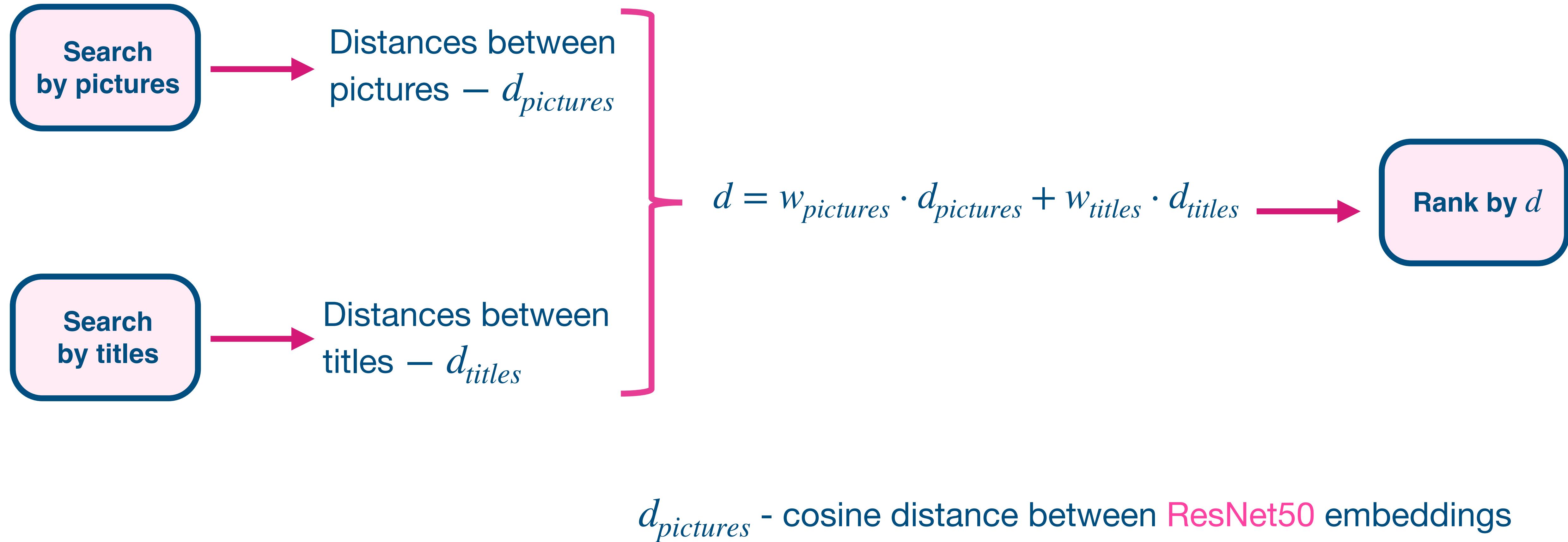
Divergence = 2.72

Divergence = 5.91

Candidates selection stages

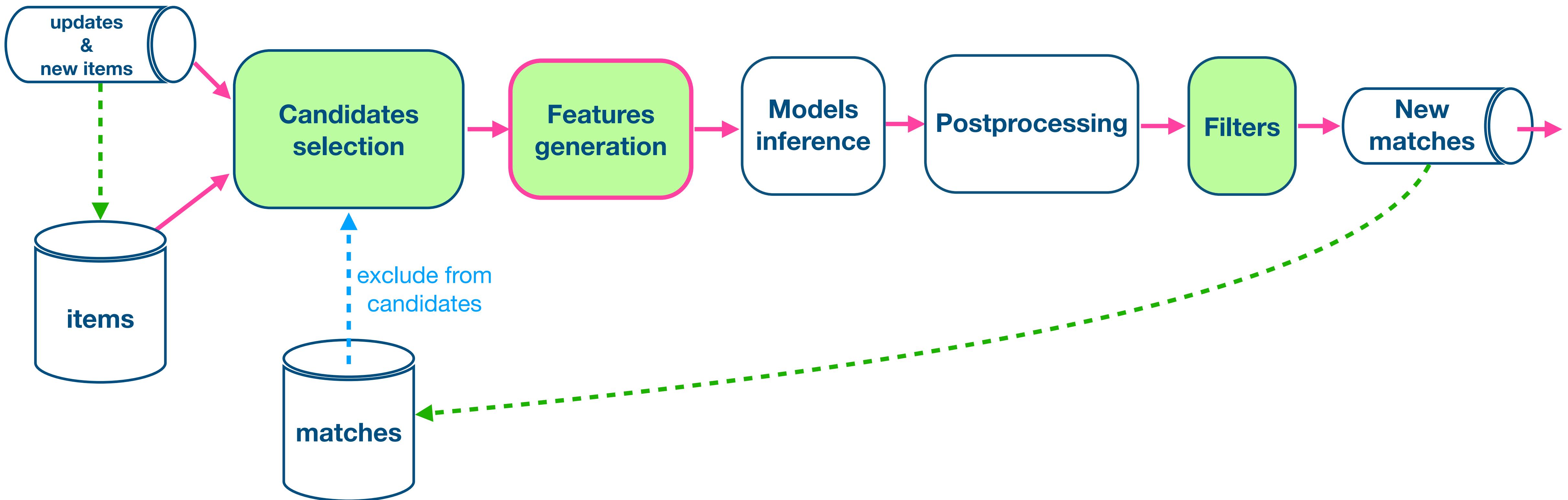


Ranking



Matching Pipeline

High level design



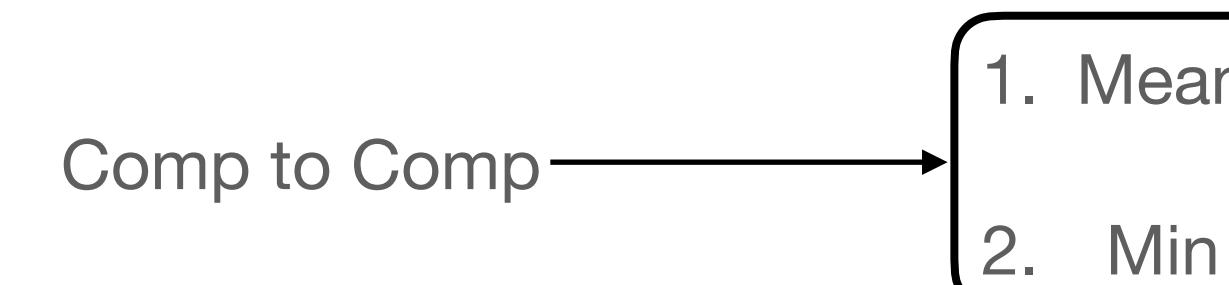
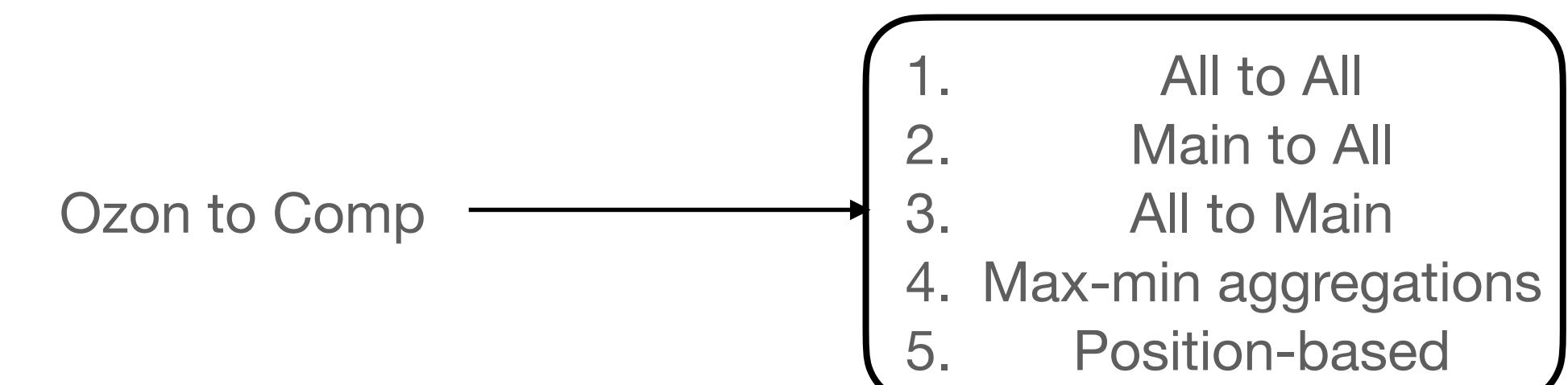
Features

Embeddings based

Ozon - Competitor matrix

Pairwise distance matrix						
	5	7	4	1	3	2
	6	6	5	4	5	6
	0	2	4	5	4	2
	2	0	6	7	5	3
	2	3	1	2	1	0
	3	6	4	5	4	4

Statistics



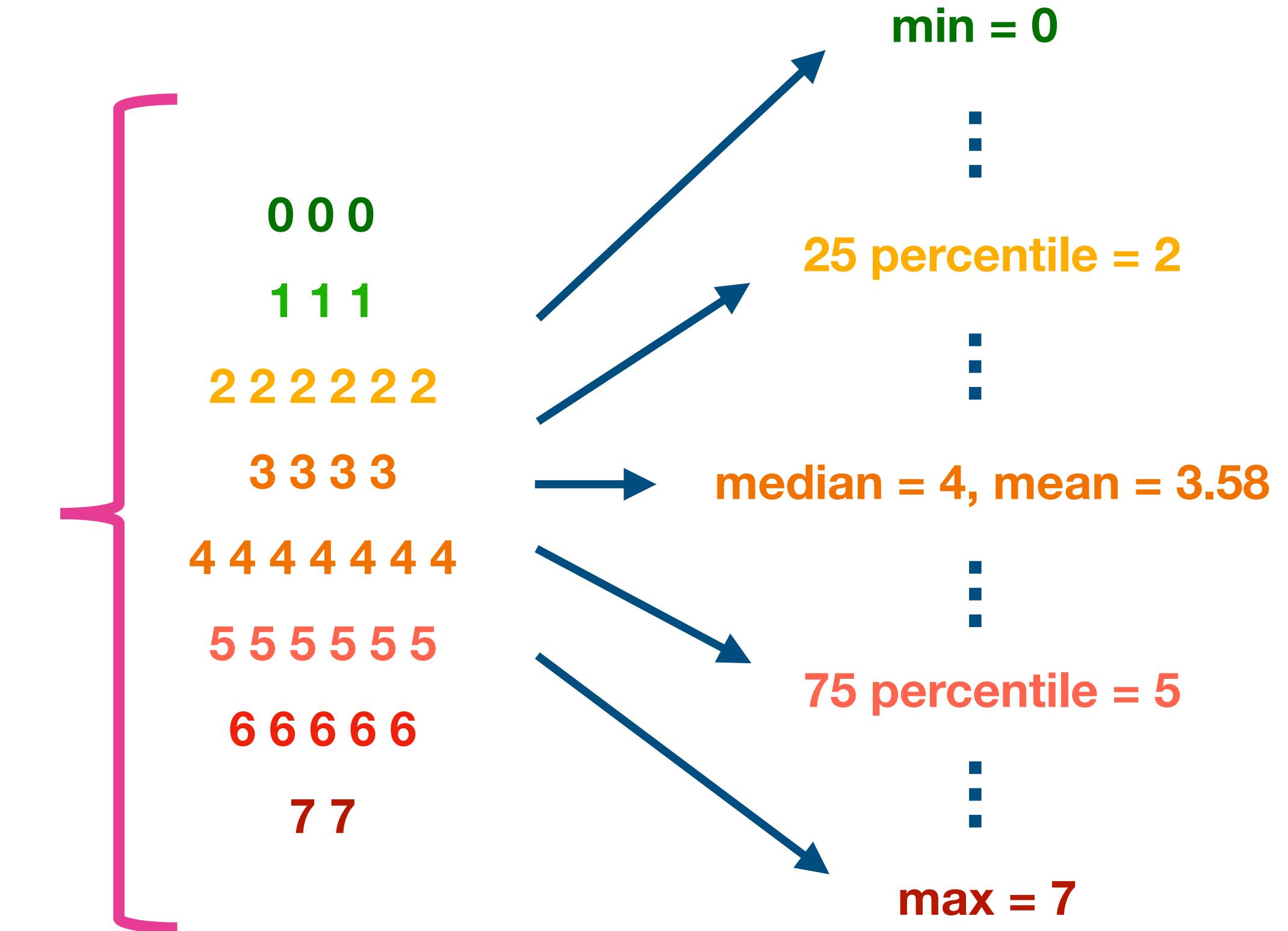
Features

Embeddings based

Ozon - Competitor matrix

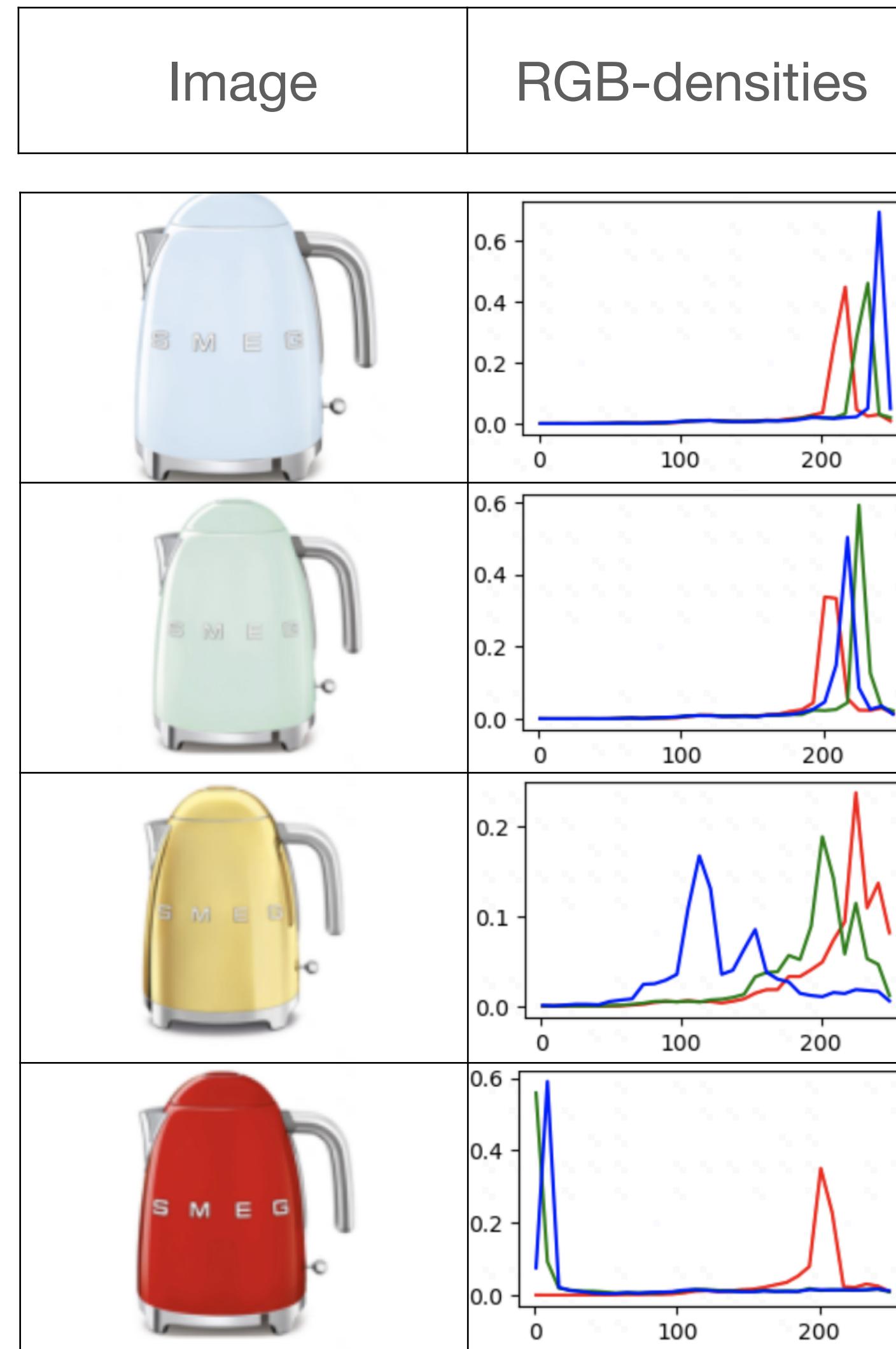
Pairwise distance matrix						
	5	7	4	1	3	2
	6	6	5	4	5	6
	0	2	4	5	4	2
	2	0	6	7	5	3
	2	3	1	2	1	0
	3	6	4	5	4	4

Statistics



Features

RGB-densities based



Jeffreys divergence:

$$D_{\text{KL}}(P \parallel Q) + D_{\text{KL}}(Q \parallel P)$$

Hellinger distance:

$$\frac{1}{\sqrt{2}} \|\sqrt{P} - \sqrt{Q}\|_2.$$

L2-norm of diff:

$$\|P - Q\|_2$$

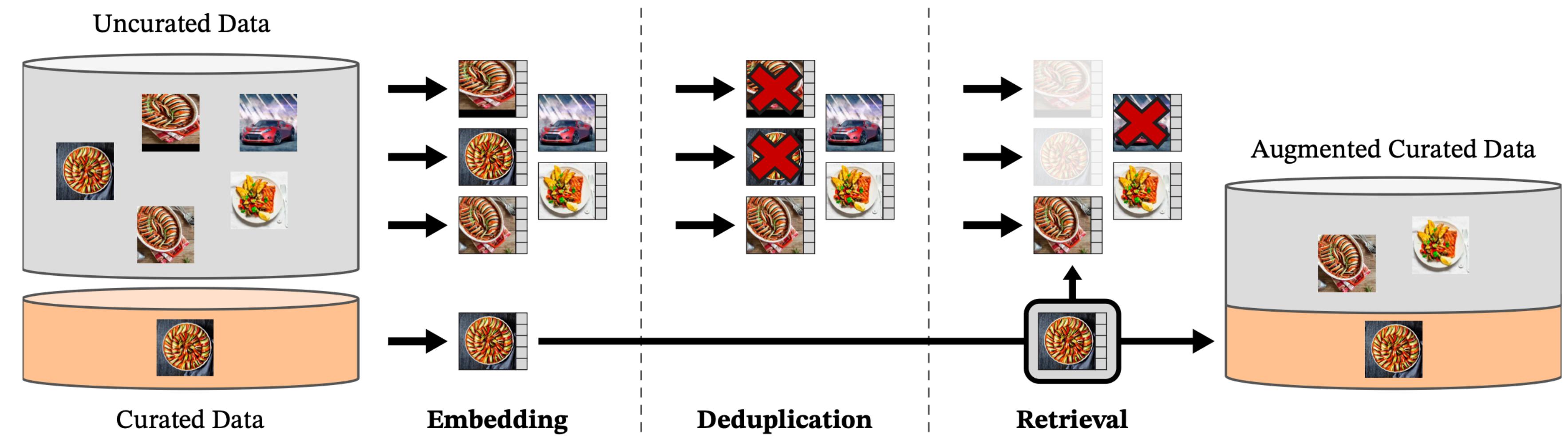
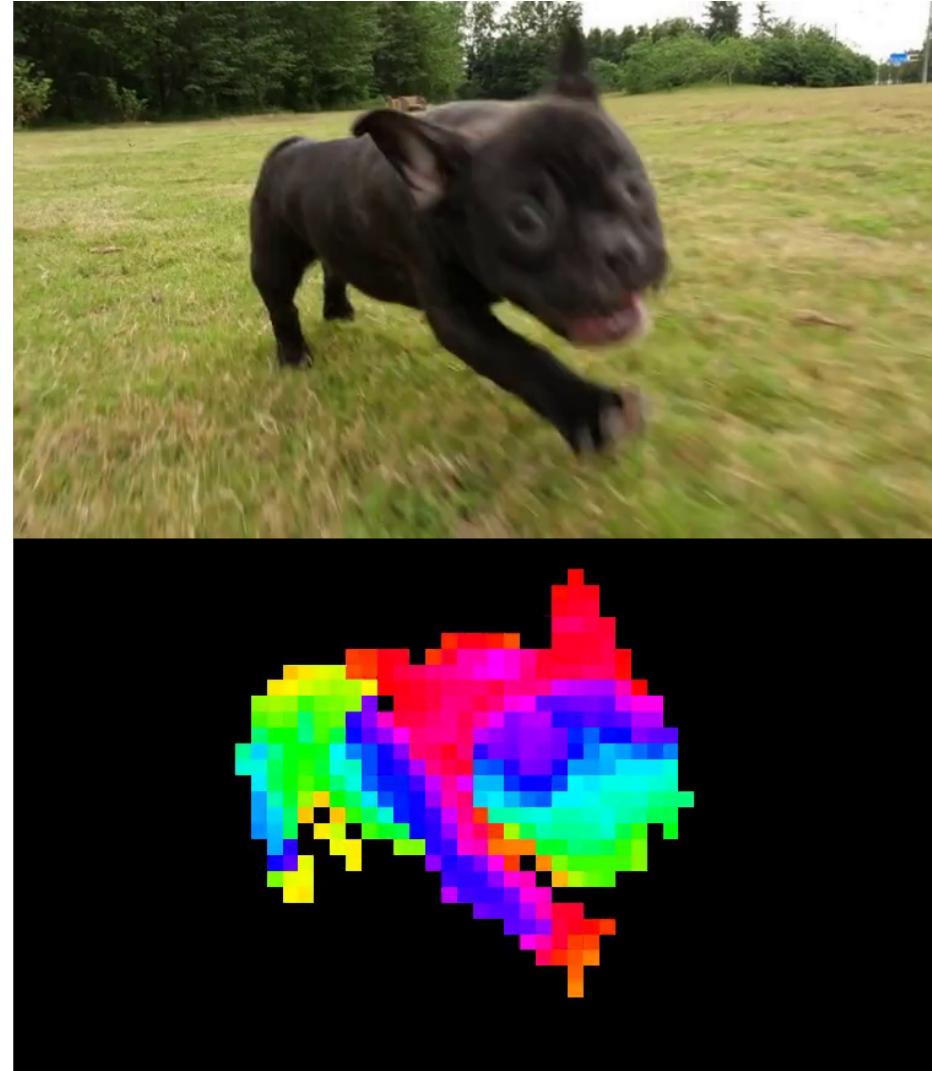
$$i_P = \text{argmax}(P), \quad i_Q = \text{argmax}(Q)$$

Weighted distance between peaks:

$$\frac{\max(P[i_P], Q[i_Q])}{\min(P[i_P], Q[i_Q])} * \frac{|i_P - i_Q|}{\text{bins}(P)}$$

Neural Networks

Current research



DINOv2

[M. Oquab et al., 2023]

Neural Networks

Current research

Print On Demand



Neural Networks

Current research

Background + info



Thanks for your attention