

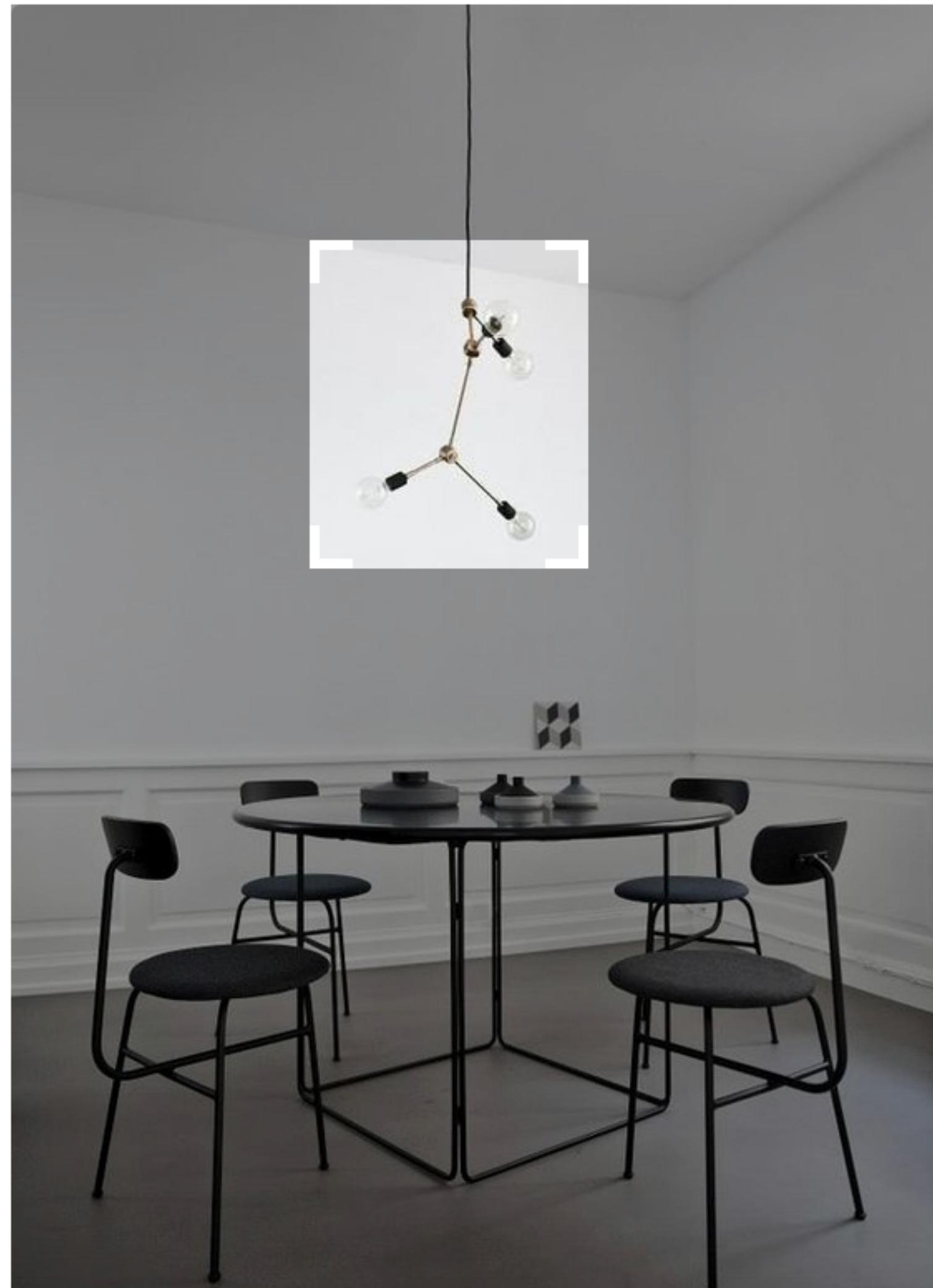


Visual Search in the Deep Learning Era

Dmitry Kislyuk
Visual Search Lead @ Pinterest

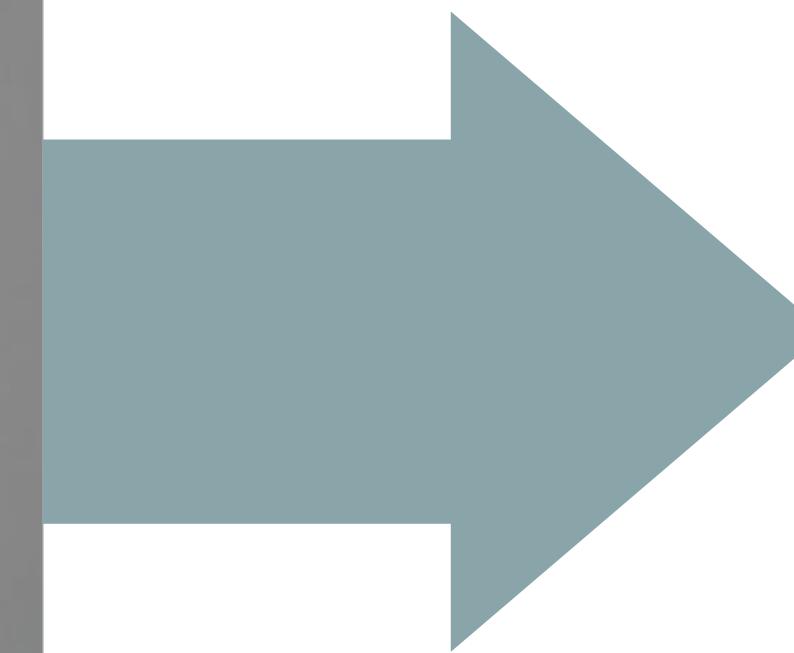
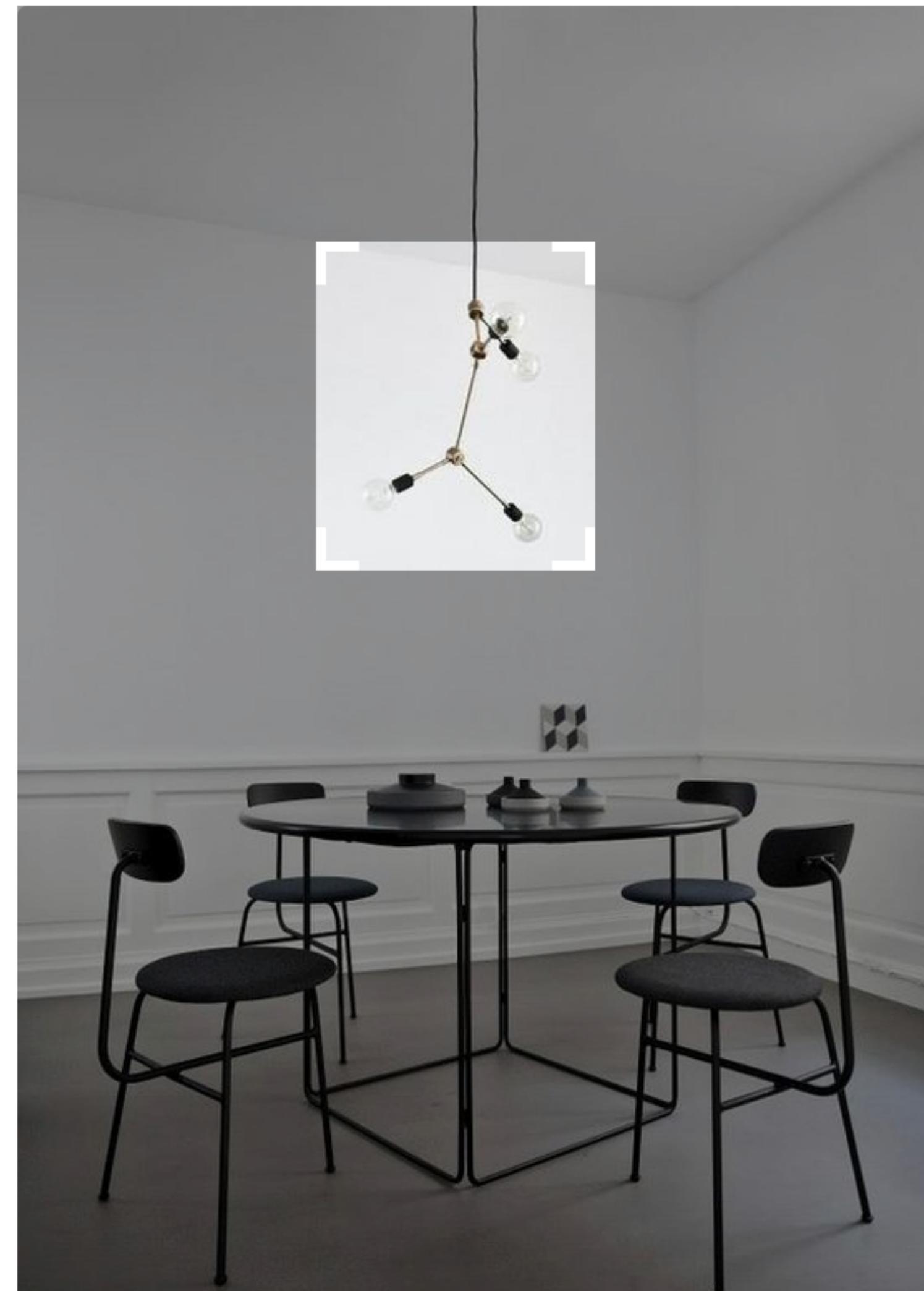
Visual Search

Visually similar results



Visual Search

Visually similar results

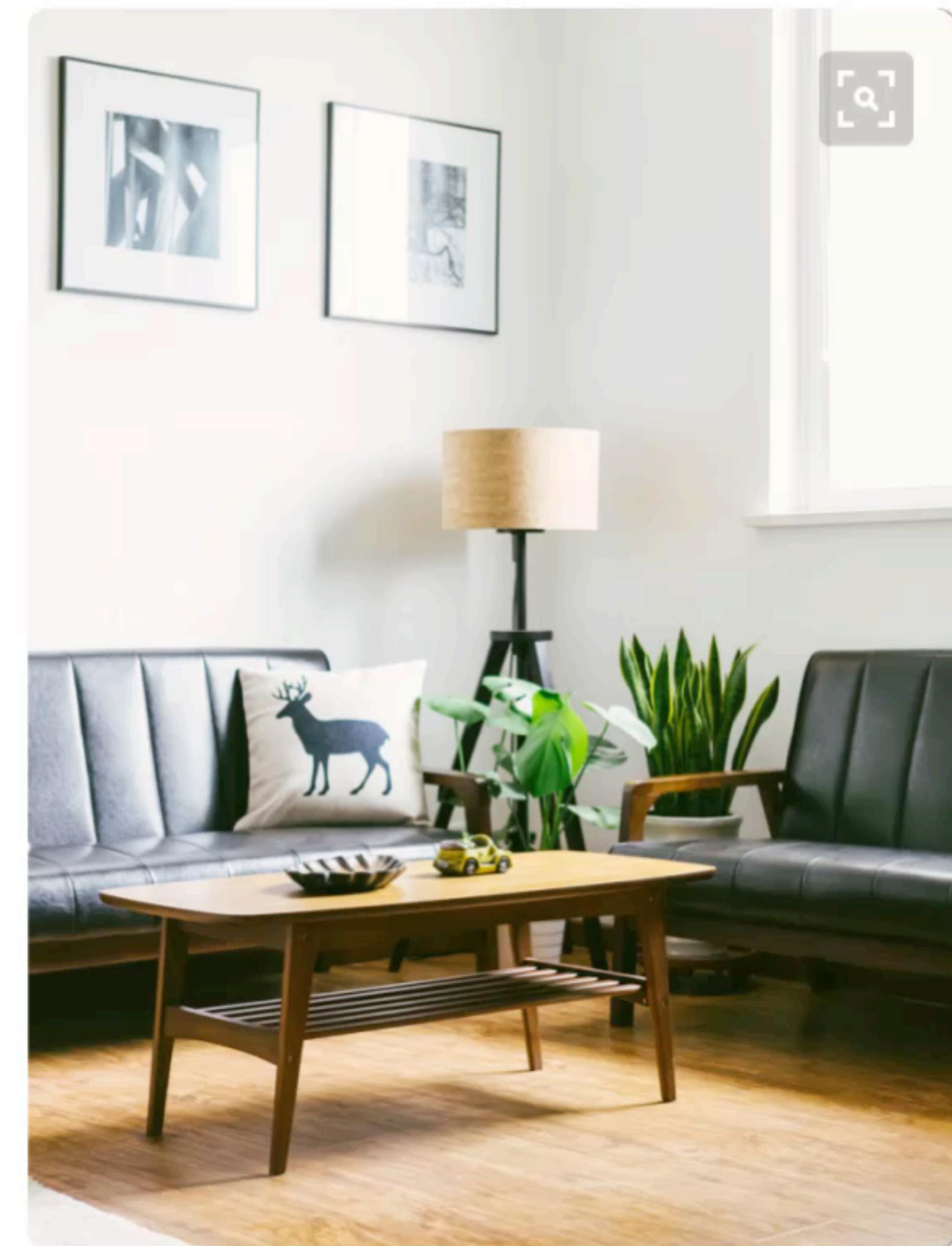


franklin chandelier tribeca franklin chandeliers franklin tribeca

Visual search results:

- New Decor
- Robin Lane Lighting
- Tamsin Paola Lighting
- from A-TAK DESIGN
A-TAK DESIGN
kasia morstyn lamps
- Veronika W Lighting
- WOHNWUNDERBAR • LAMPEN •
- Carla Vent Home
- Karin Daar Belysning
- Heidi Risku
- Magnificently Modern
Modern & Contemporary ...
- Nest.co.uk

Visual Search



Uploaded by
Ben Chiaramonte

Read it

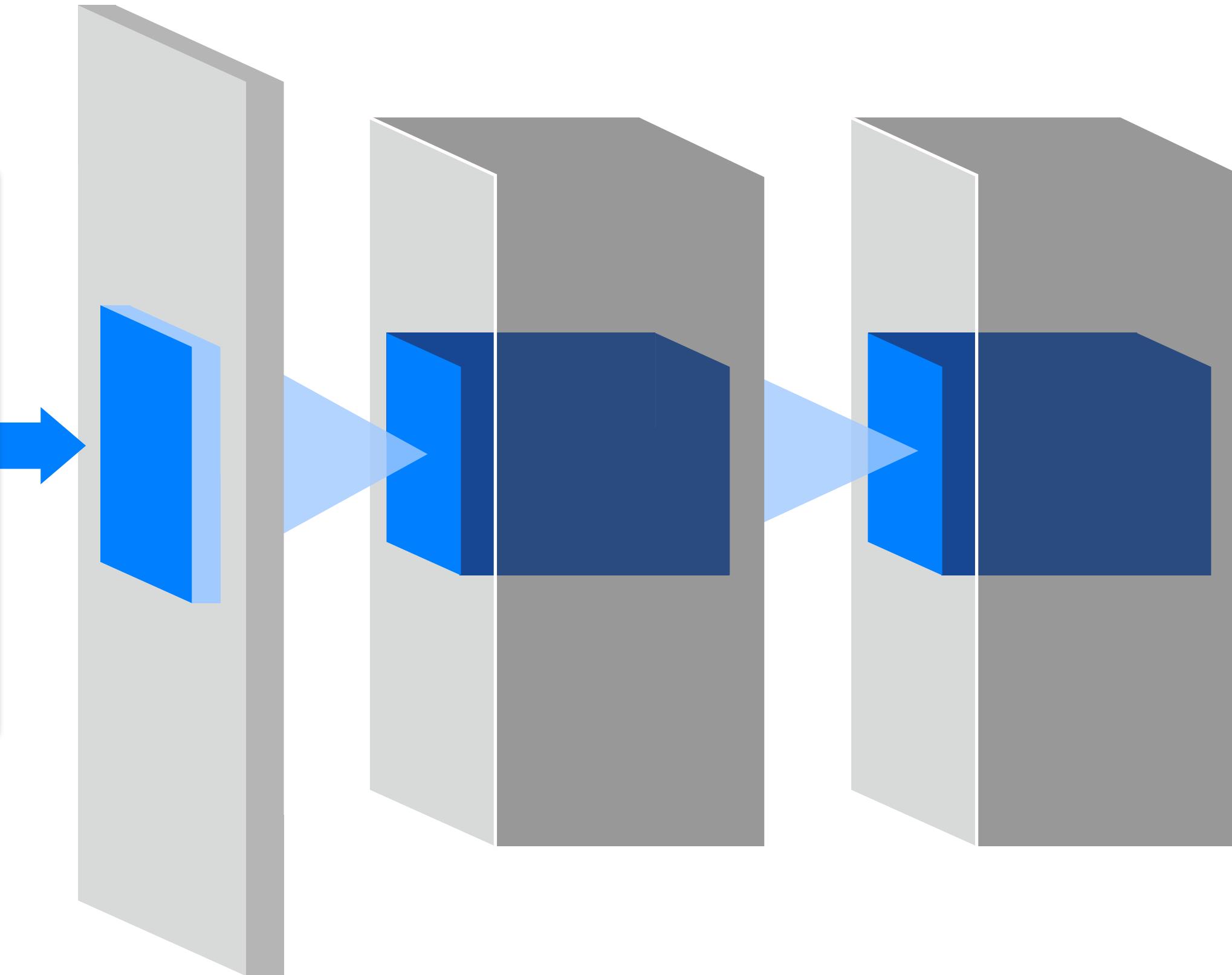
Add a comment



Technical Challenge

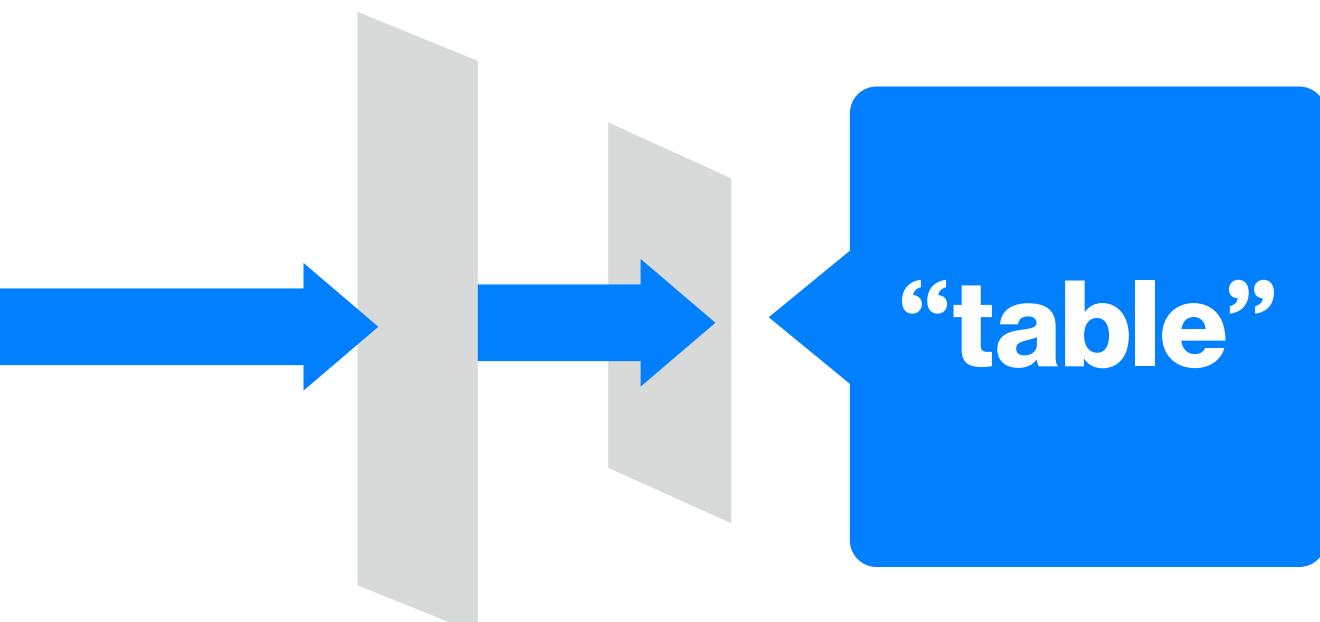
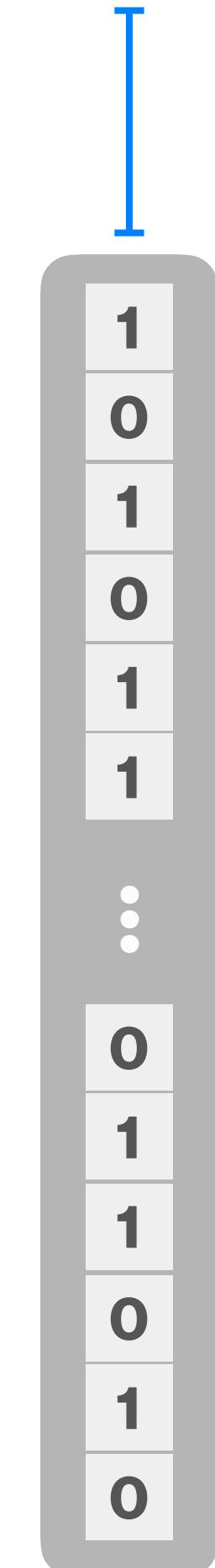
Define visual similarity
between any visual
object and images in a
dataset, in real time.

Convolutional networks for image classification



Extract high level features

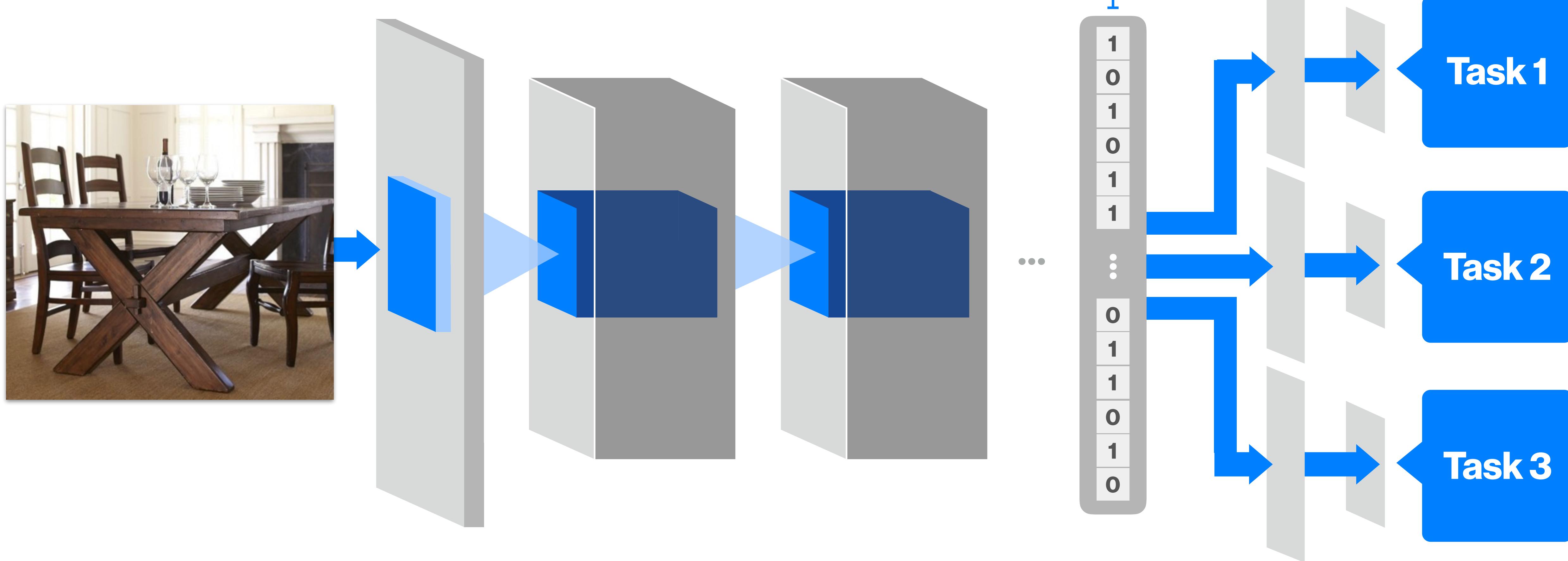
high-dimensional embedding



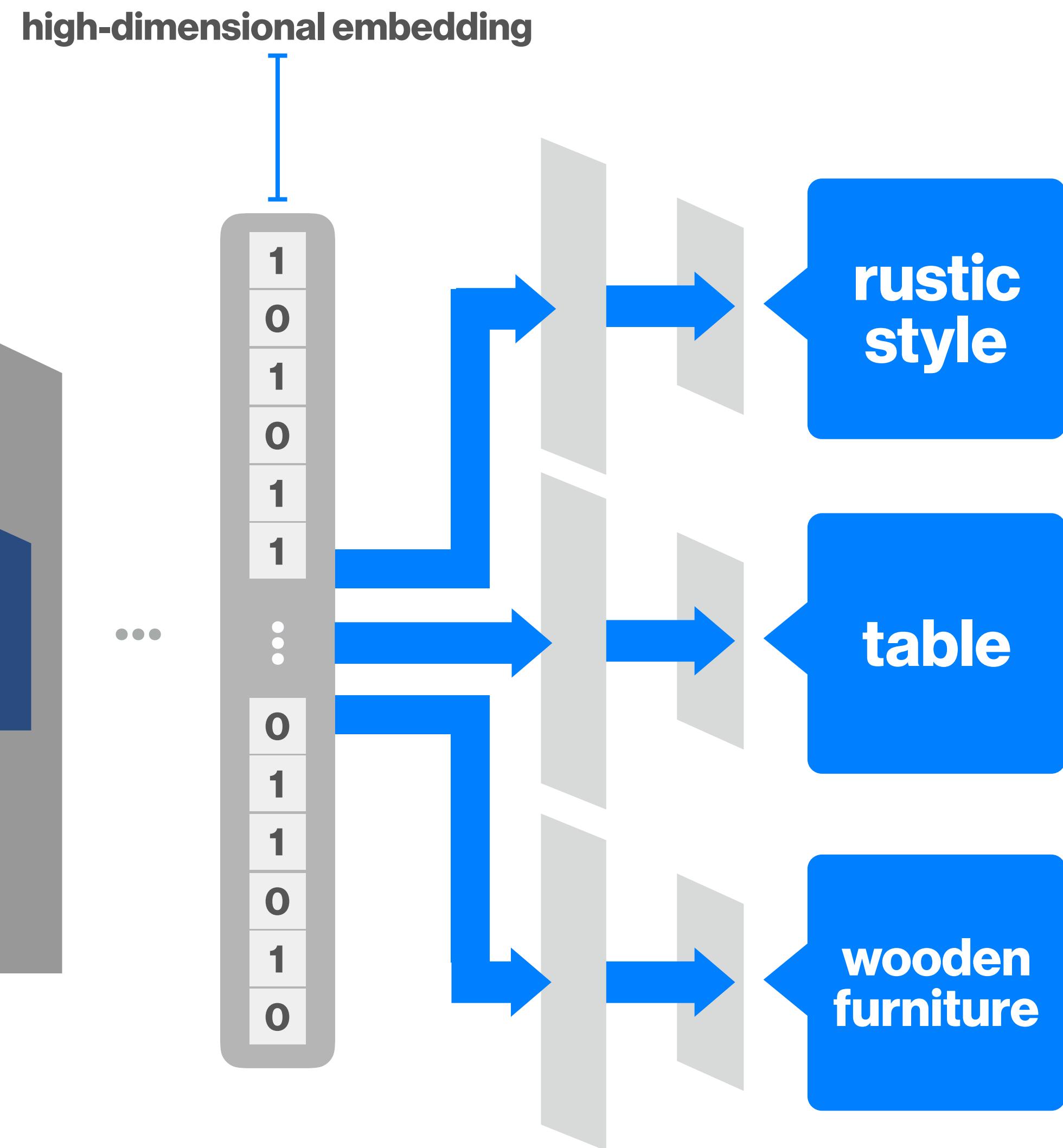
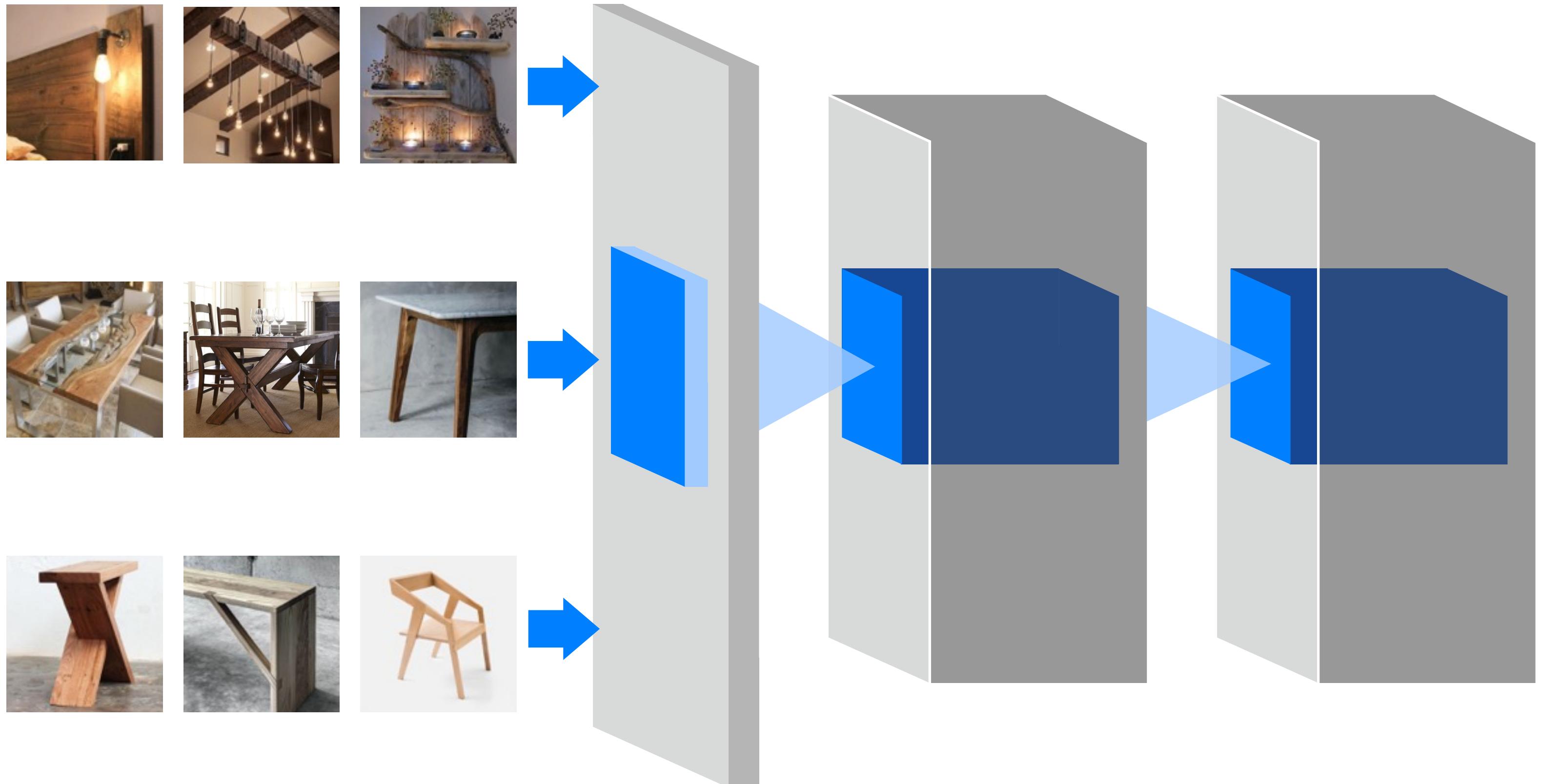
Representation

Classify Each Sample

Convolutional networks for multi-task learning



Convolutional networks for multi-task learning



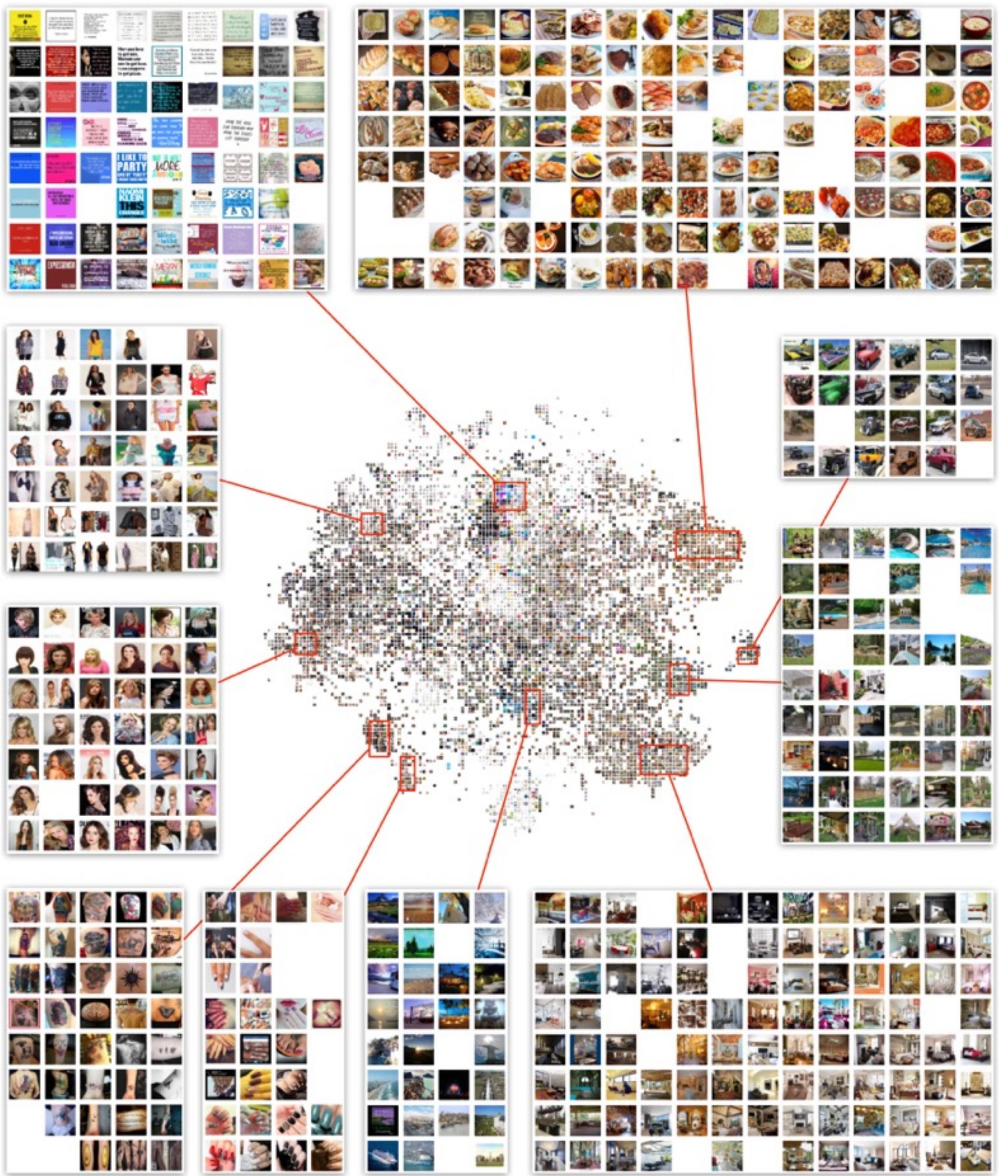
Extract high level features

Representation

Multi-task Learning

Deep visual features

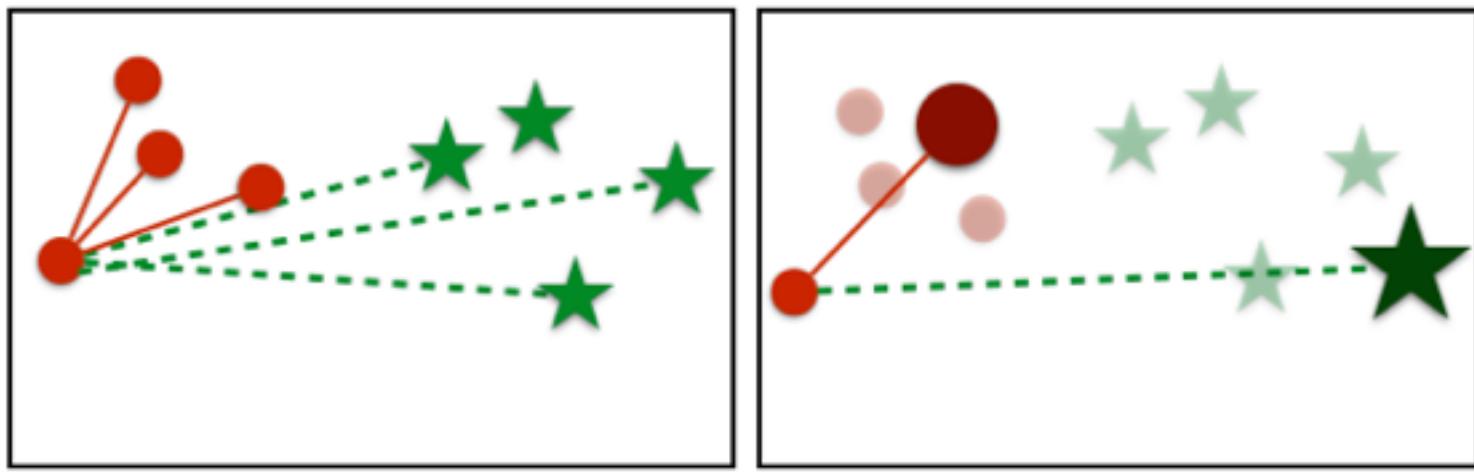
t-SNE of fine-tuned, binarized VGG-16 features



Choice of embedding model

Kislyuk et al., WWW'17

Proxy Ranking Loss (Movshovitz-Attias et al.) + Multi-Task



$$L_{\text{triplet}}(x, y, z) = [d(x, y) + M - d(x, z)]_+$$

Other embedding models

ResNet-152 / ResNeXt embeddings lead to best overall p@k performance, but features 4x larger

Binarized **VGG-16 fc6** most compact high-performing embeddings

Fine-tuning on Pinterest data helps further

Model	Layer	Type	Dist.	P@1	P@5	P@10
AlexNet	fc6	raw	L2	0.093	0.045	0.027
AlexNet	fc6	raw	L1	0.099	0.047	0.027
AlexNet	fc7	raw	L2	0.104	0.052	0.031
AlexNet	fc7	raw	L1	0.106	0.052	0.031
AlexNet	fc8	raw	L2	0.088	0.047	0.028
AlexNet	fc8	raw	L1	0.090	0.046	0.028
GoogleNet	loss3/classifier	raw	L2	0.095	0.050	0.032
GoogleNet	loss3/classifier	raw	L1	0.098	0.050	0.032
VGG16	fc6	raw	L2	0.108	0.051	0.030
VGG16	fc6	raw	L1	0.118	0.057	0.035
VGG16	fc7	raw	L2	0.116	0.058	0.036
VGG16	fc7	raw	L1	0.113	0.060	0.038
VGG16	fc8	raw	L2	0.104	0.054	0.034
VGG16	fc8	raw	L1	0.106	0.054	0.034
ResNet101	pool5	raw	L2	0.160	0.080	0.050
ResNet101	pool5	raw	L1	0.149	0.073	0.045
ResNet101	fc1000	raw	L2	0.133	0.068	0.042
ResNet101	fc1000	raw	L1	0.139	0.067	0.041
ResNet152	pool5	raw	L2	0.170	0.083	0.050
ResNet152	pool5	raw	L1	0.152	0.077	0.047
ResNet152	fc1000	raw	L2	0.149	0.073	0.045
ResNet152	fc1000	raw	L1	0.148	0.073	0.044
AlexNet	fc6	binary	H.	0.129	0.065	0.039
AlexNet	fc7	binary	H.	0.110	0.054	0.033
AlexNet	fc8	binary	H.	0.089	0.046	0.027
VGG16	fc6	binary	H.	0.158	0.081	0.049
VGG16	fc7	binary	H.	0.133	0.068	0.044
VGG16	fc8	binary	H.	0.110	0.055	0.035
ResNet101	fc1000	binary	H.	0.125	0.062	0.039
ResNet101	pool5	binary	H.	0.055	0.025	0.014
ResNet152	fc1000	binary	H.	0.133	0.065	0.041
ResNet152	pool5	binary	H.	0.057	0.026	0.015
VGG16 (Pin.)	fc6	binary	H.	0.169	0.089	0.056

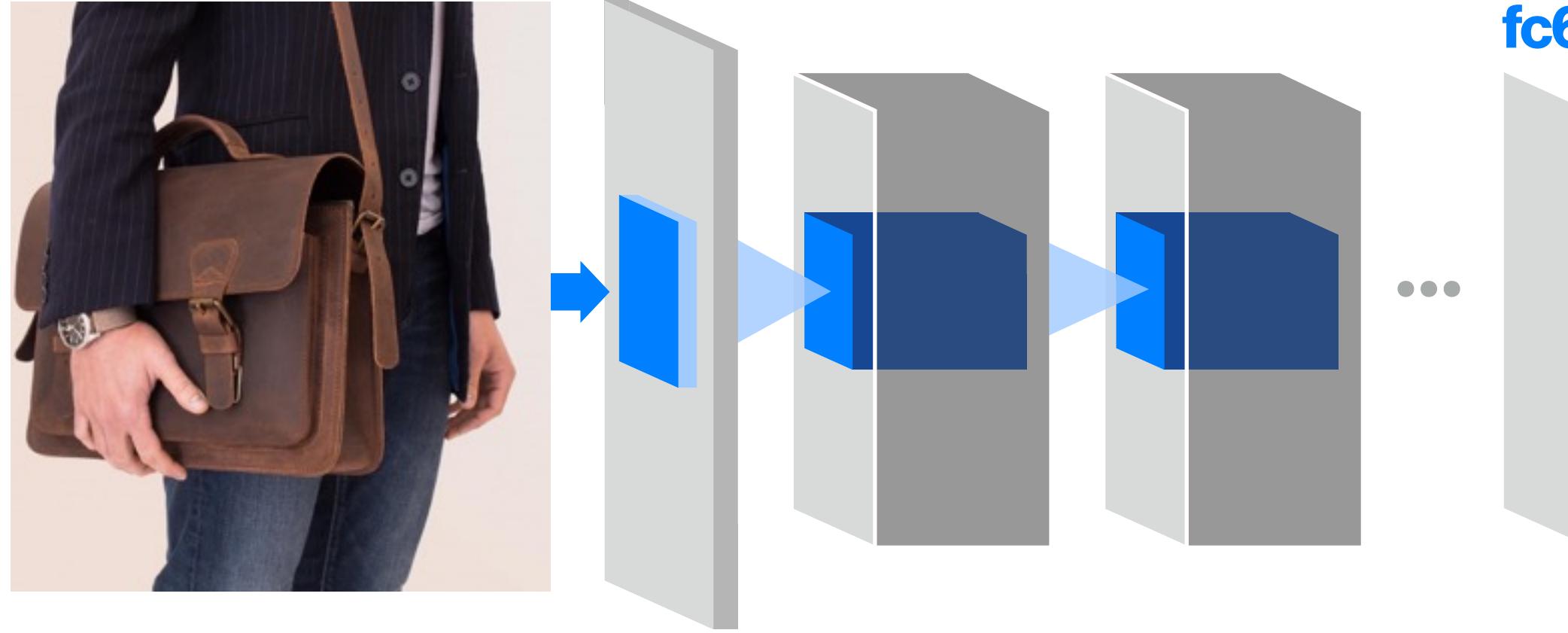
Determining visual similarity



Similarity (A,B)

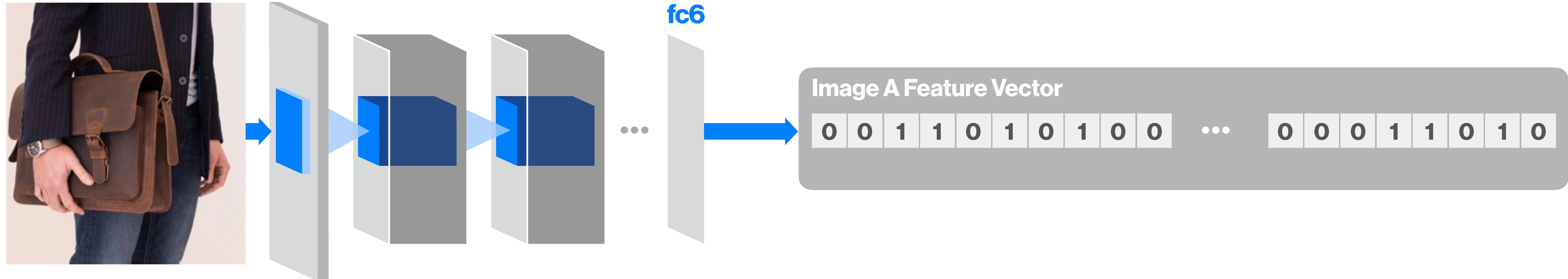


Determining visual similarity

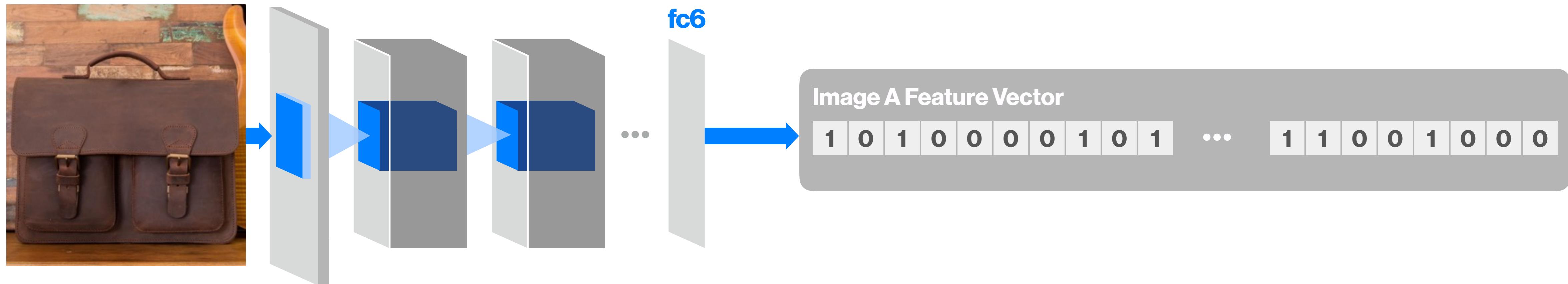
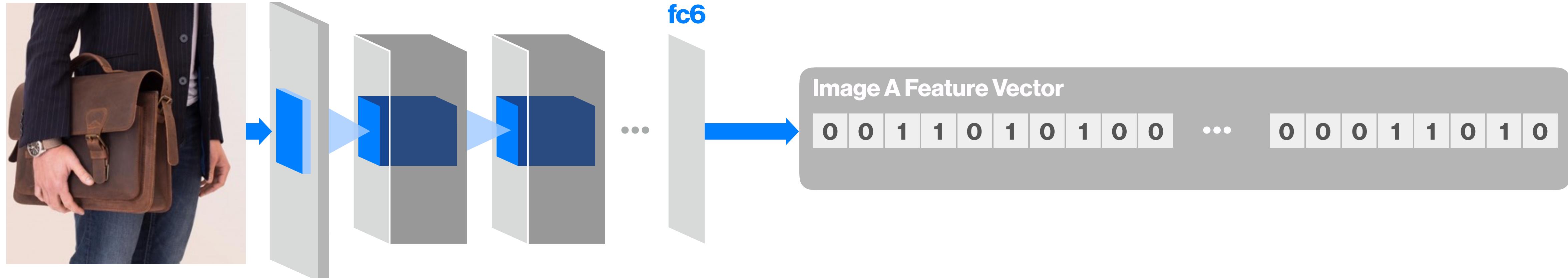


Similarity (A,B)

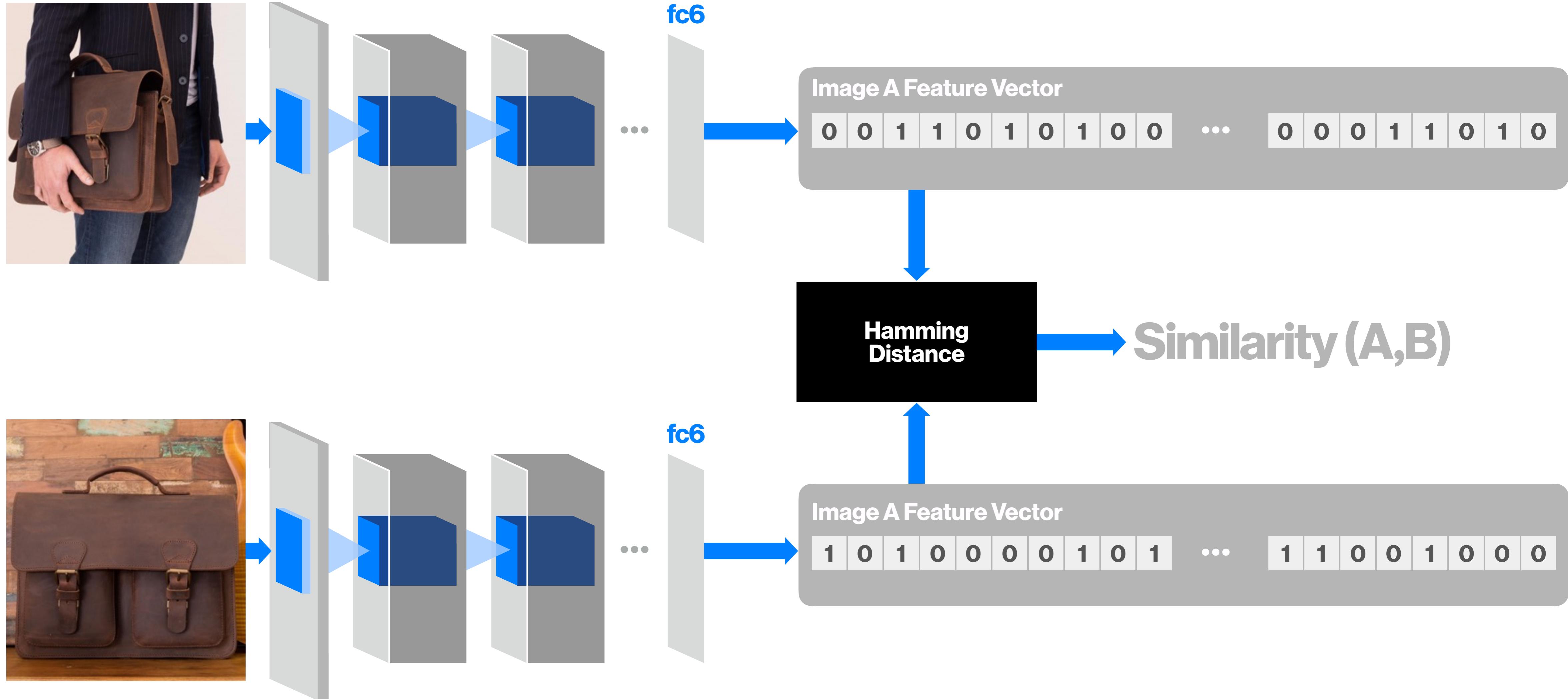
Determining visual similarity



Determining visual similarity



Determining visual similarity

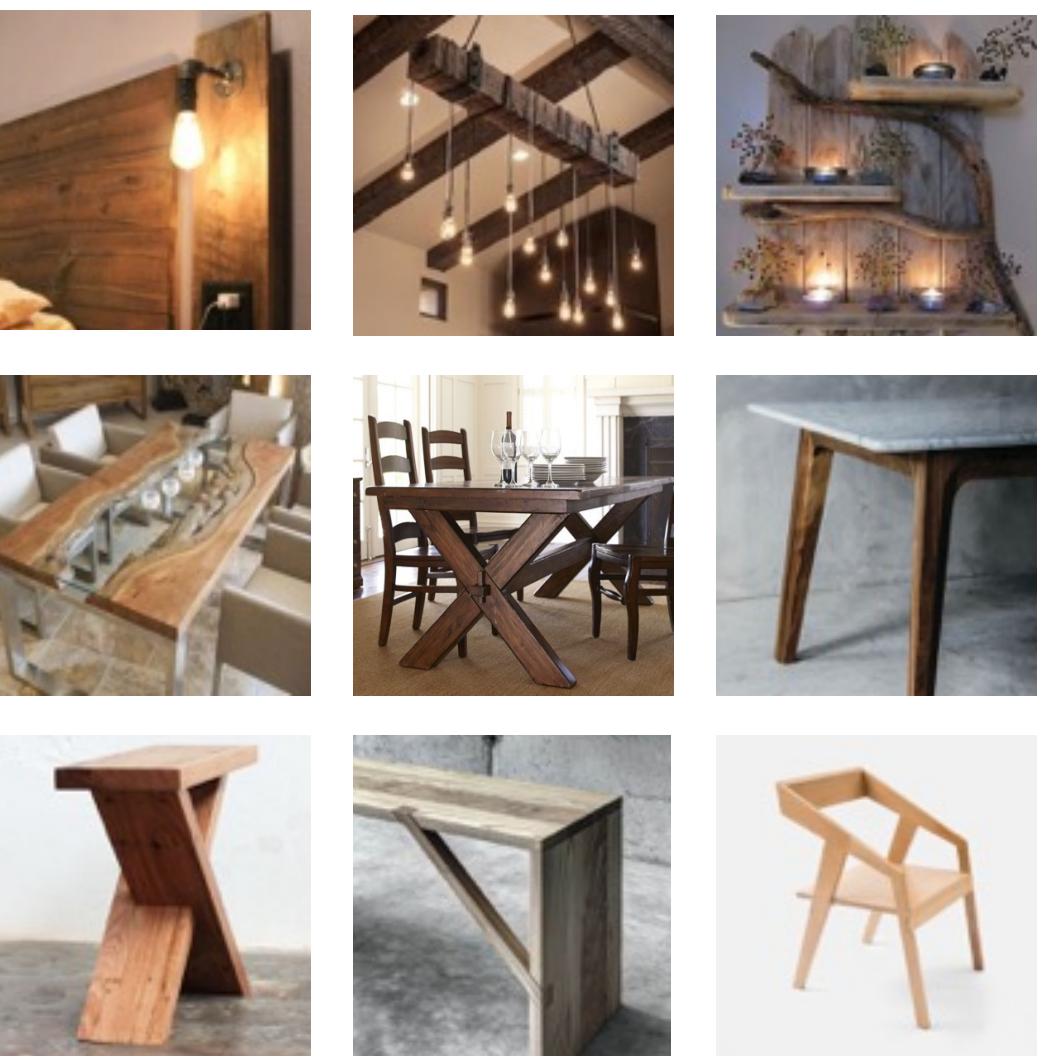


Determining visual similarity

Scaling to billions of images

- Collect training examples from different tasks (predict texture, semantic label, color, etc.)

Rustic Style



Table

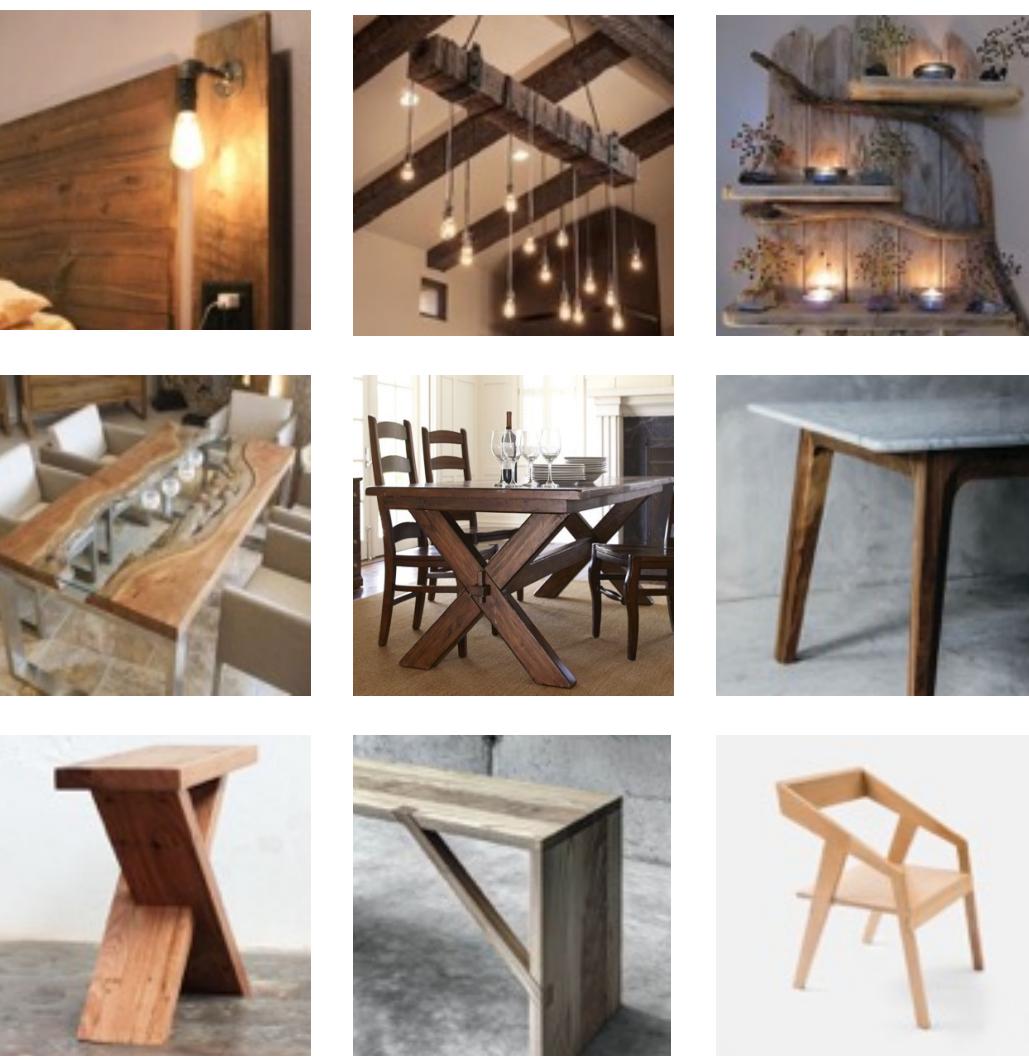
Wooden Furniture

Determining visual similarity

Scaling to billions of images

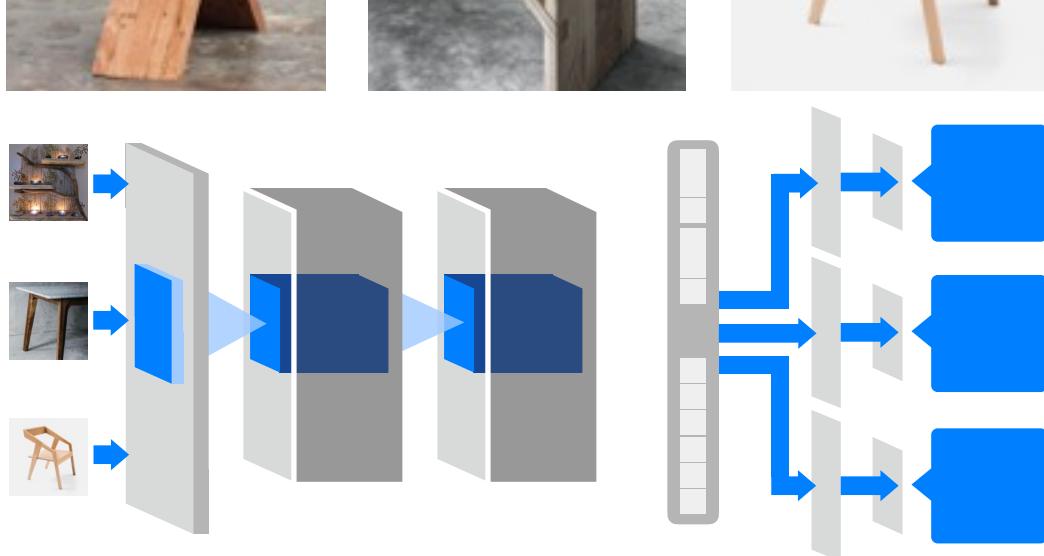
- Collect training examples from different tasks (predict texture, semantic label, color, etc.)
- Train multi-task embedding model on labeled data

Rustic Style



Table

Wooden Furniture

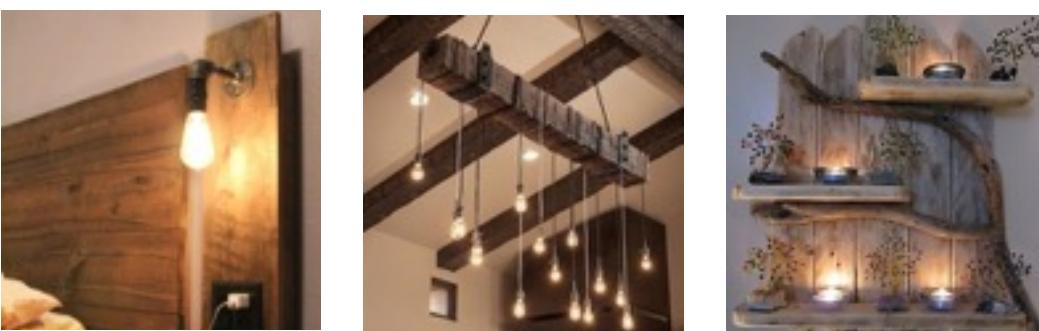


Determining visual similarity

Scaling to billions of images

- Collect training examples from different tasks (predict texture, semantic label, color, etc.)
- Train multi-task embedding model on labeled data
- Extract embedding (intermediate convolutional net features) from trained model for the **entire corpus** of images (billions)

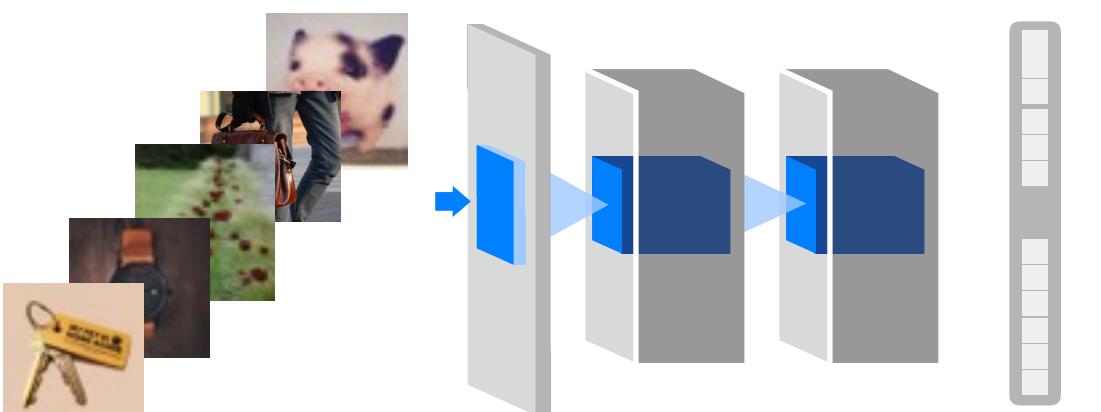
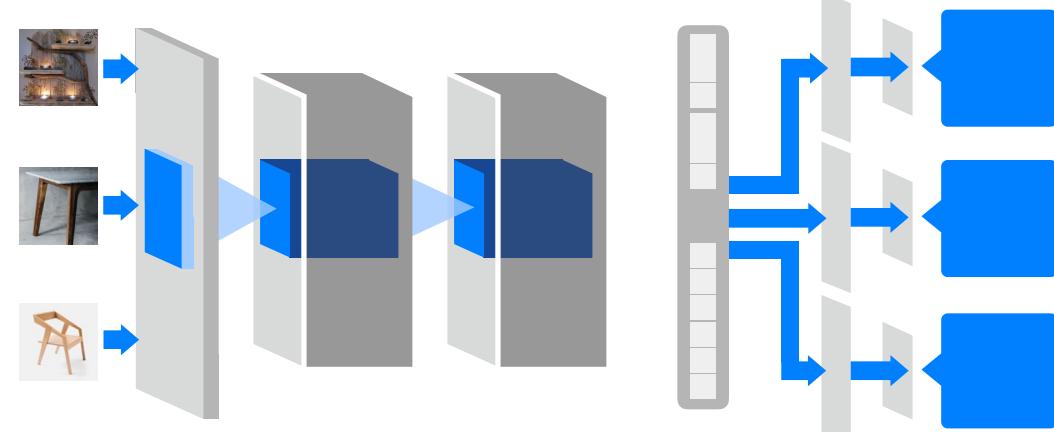
Rustic Style



Table



Wooden Furniture

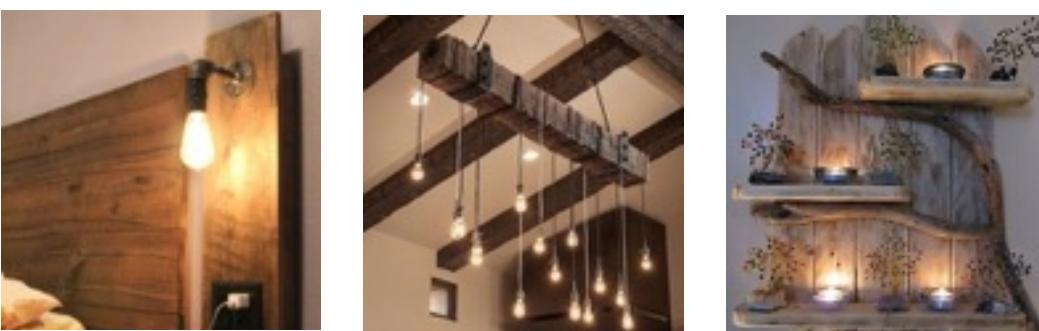


Determining visual similarity

Scaling to billions of images

- Collect training examples from different tasks (predict texture, semantic label, color, etc.)
- Train multi-task embedding model on labeled data
- Extract embedding (intermediate convolutional net features) from trained model for the **entire corpus** of images (billions)
- Index image embeddings into Approximate Nearest Neighbor (ANN) data structure (distributed!)

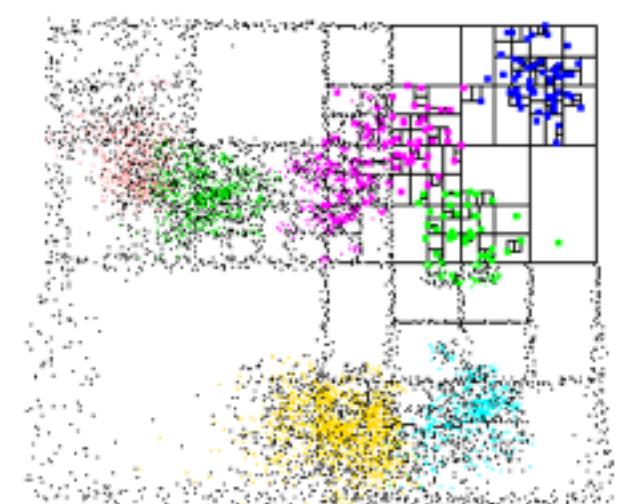
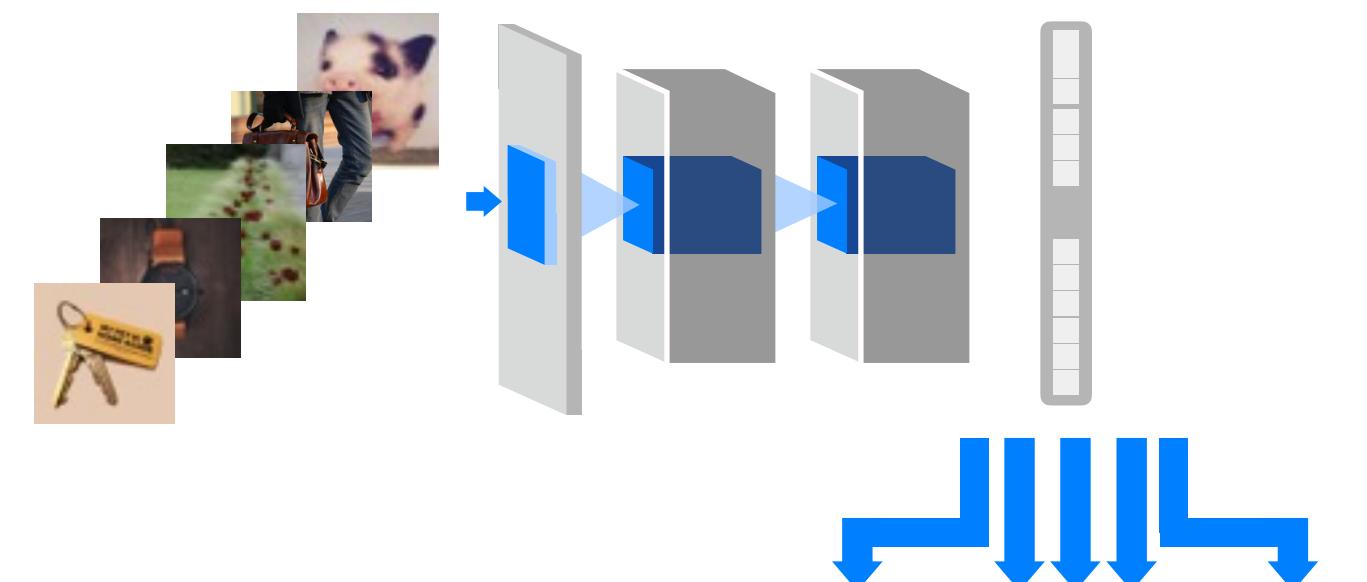
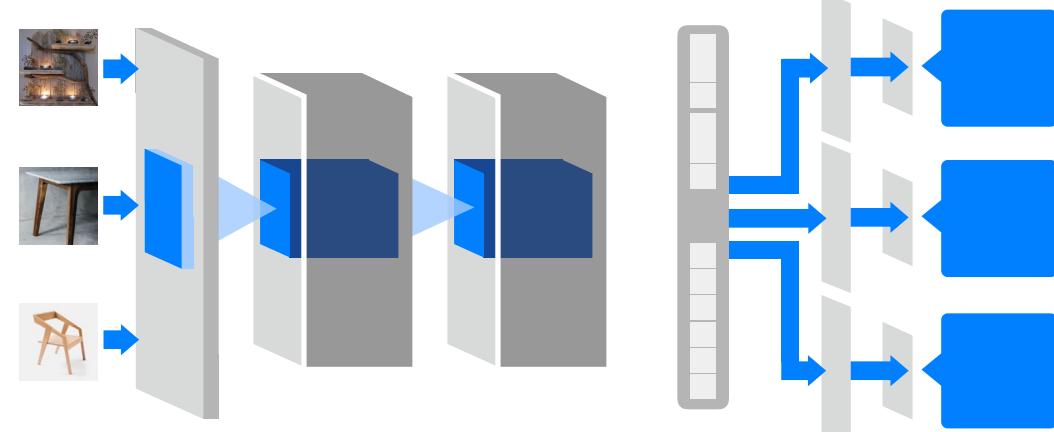
Rustic Style



Table

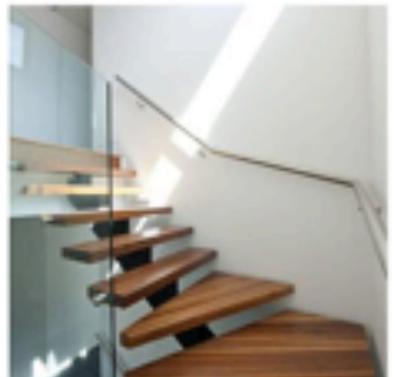


Wooden Furniture

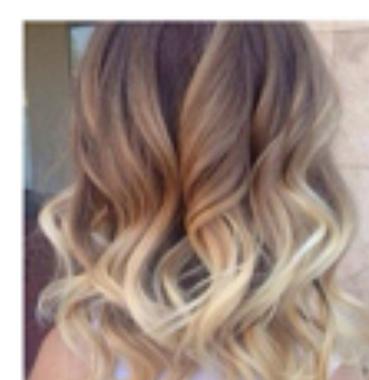


Visual Search on Pinterest with distributed nearest neighbors

Query

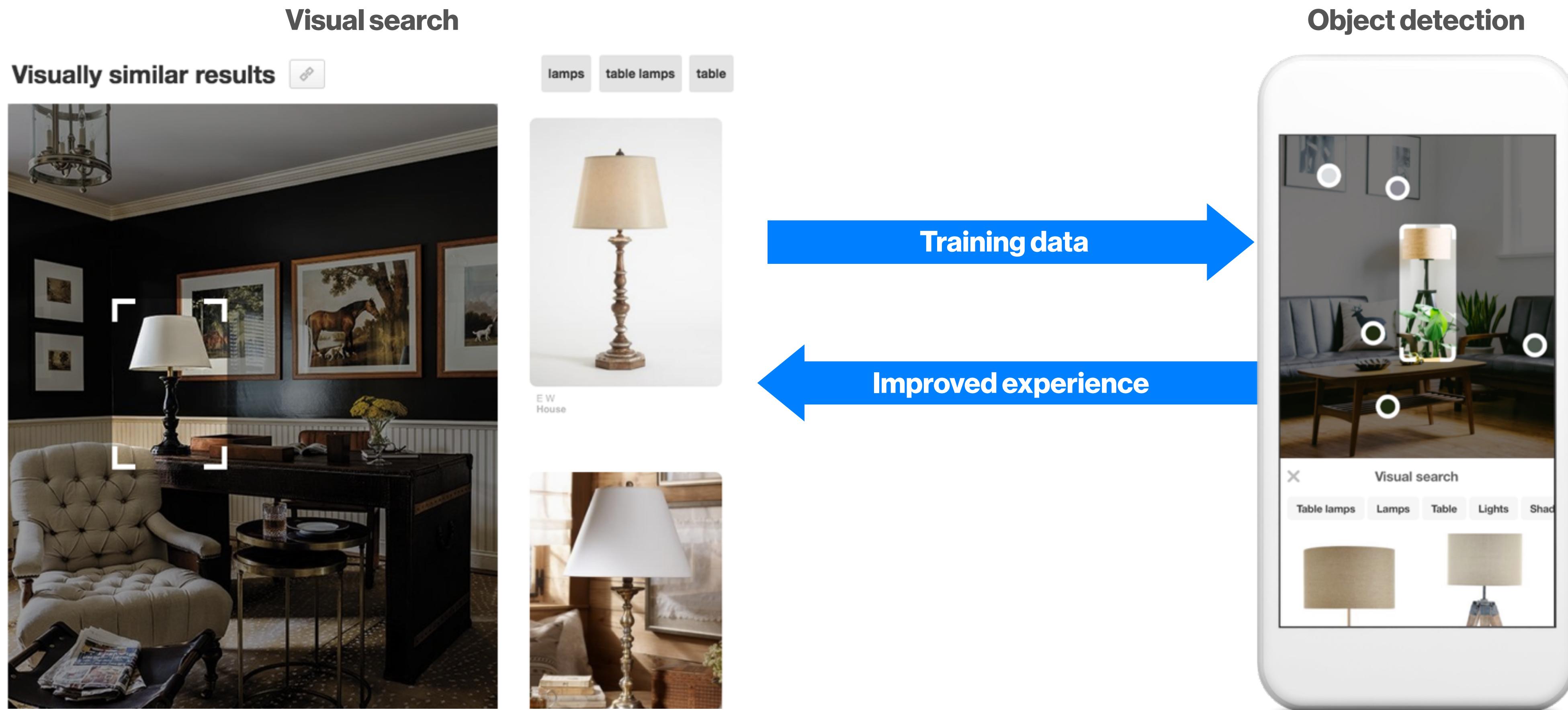


Query

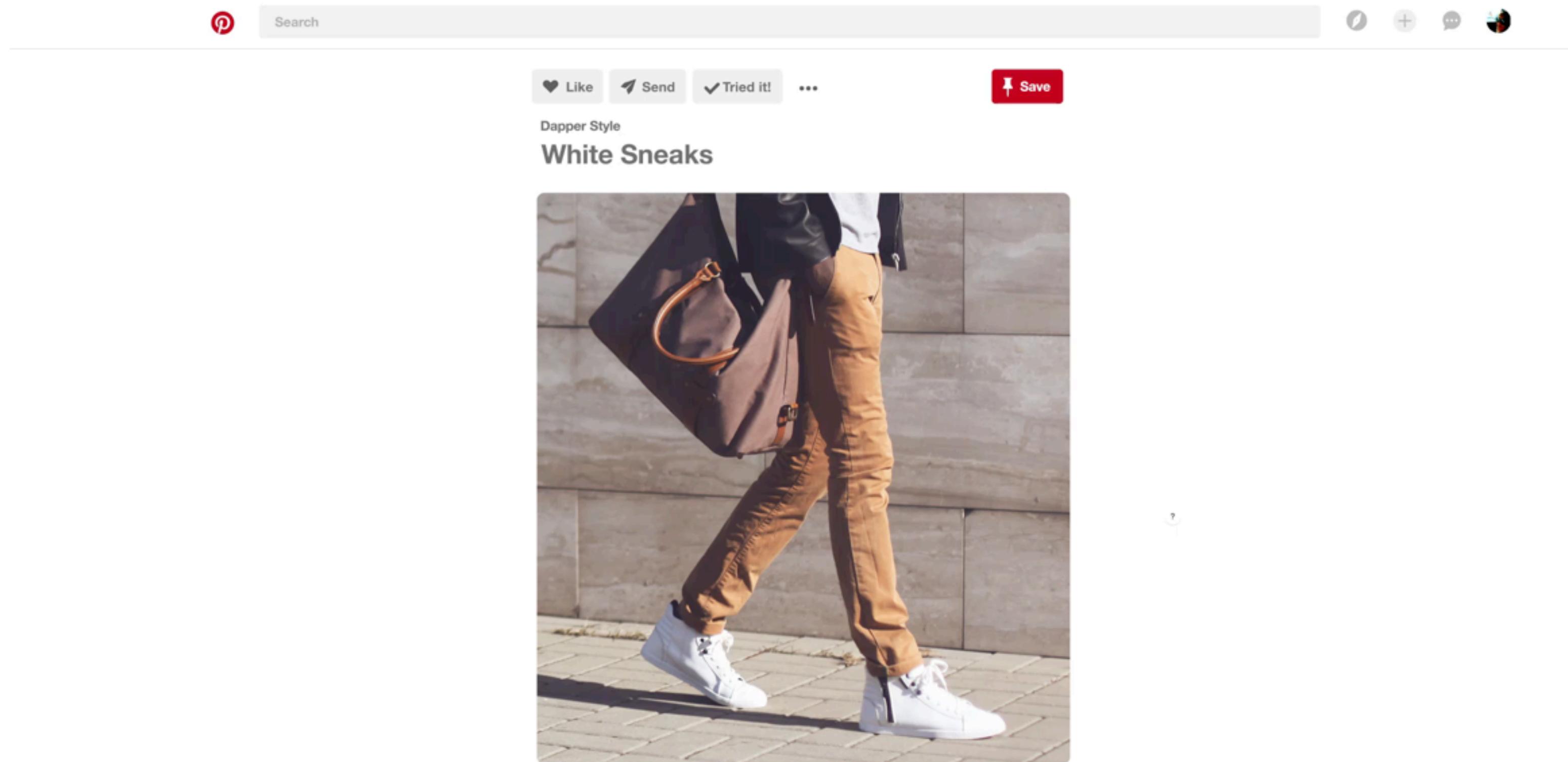


Visual Search Ecosystem

leveraging object detection

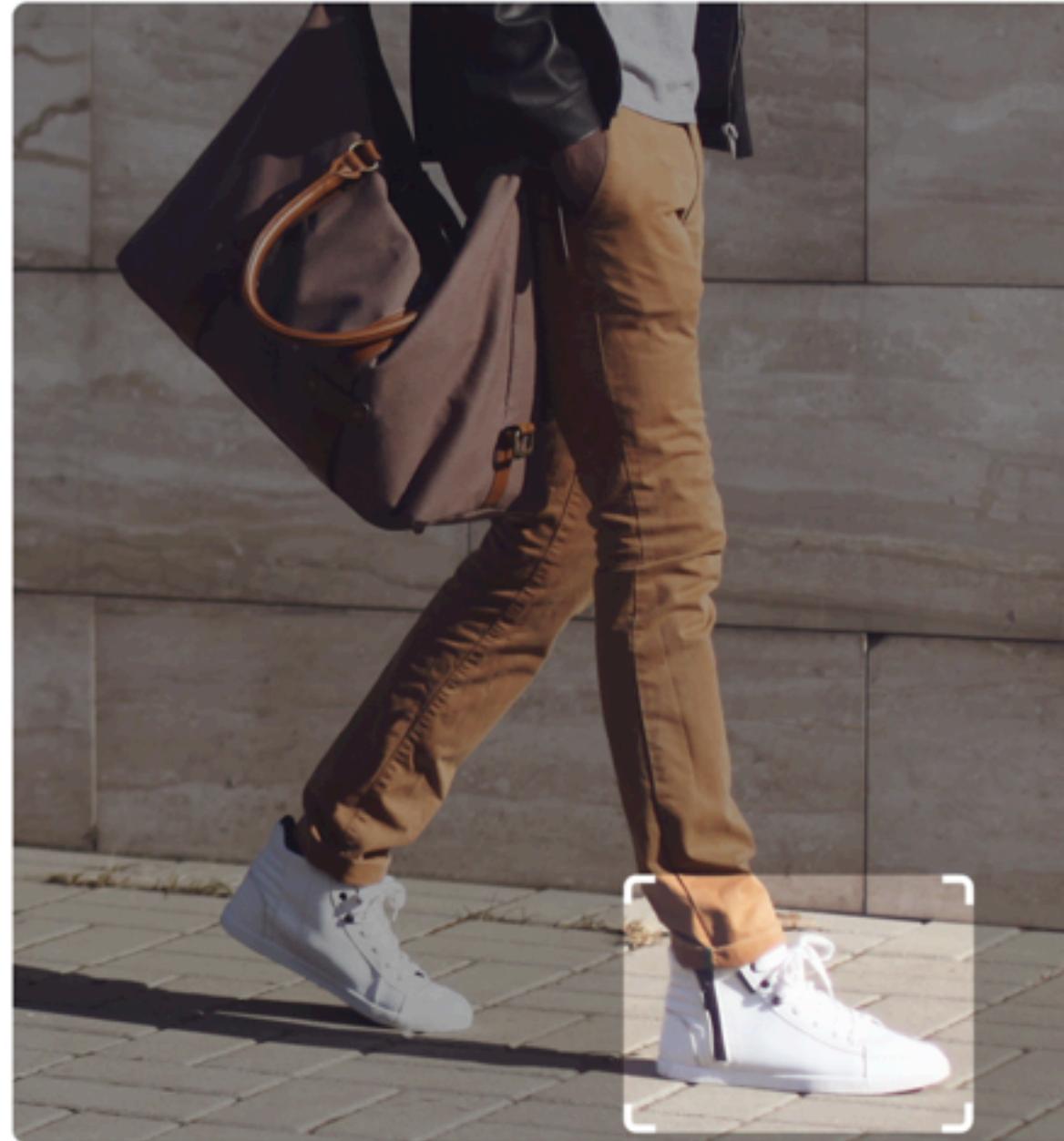


Improving detection with search



Improving detection with search

Visually similar results

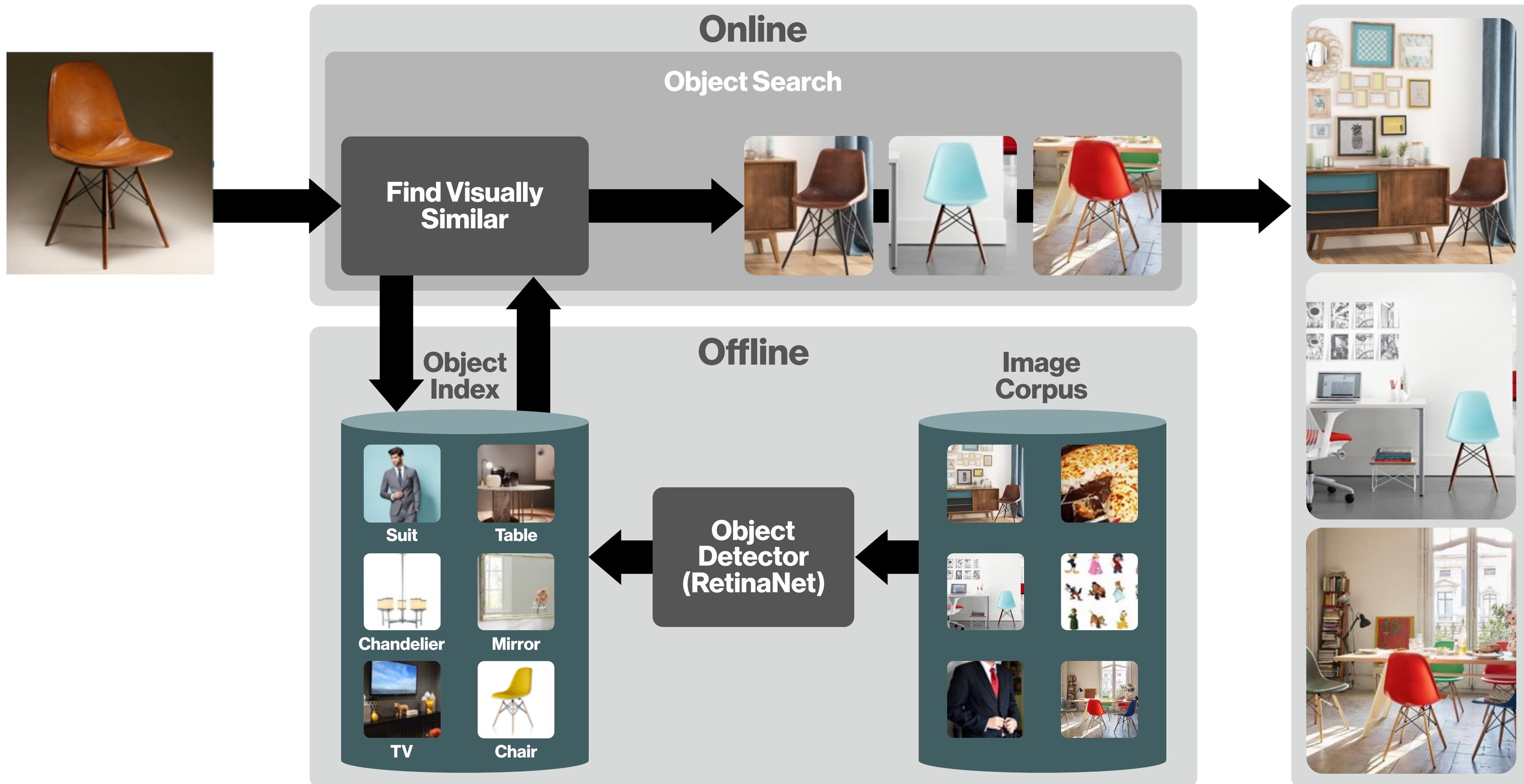


shoes sneakers nike adidas fashion light up shoes style air force

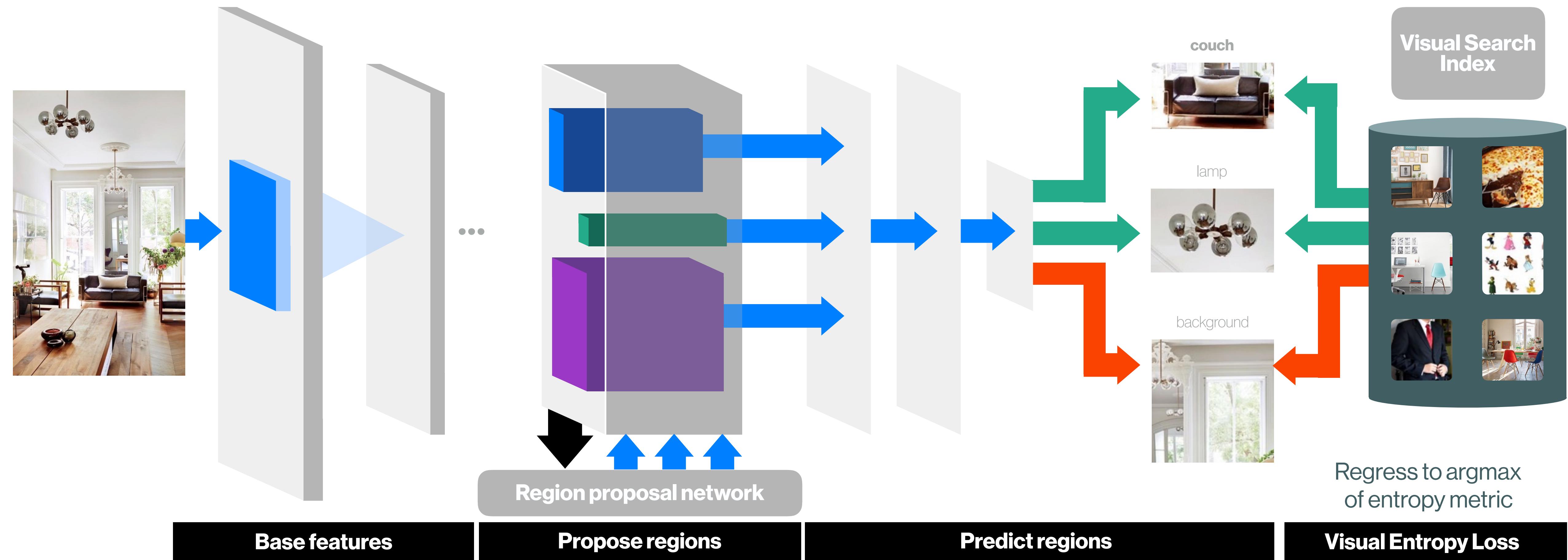
A search interface displaying a grid of visually similar shoe images. The interface includes a header with tags: shoes, sneakers, nike, adidas, fashion, light up shoes, style, and air force. Below the header, there are two rows of three images each. Each image shows a person's legs and feet wearing white sneakers. The images are accompanied by captions and user information:

- Nike**: Kasia fashion
- V**: Gabriela Sg #15
- "zizi repetto"**: Bonnie & Jane Look
- kris van assche sneakers**: Natalia Biliska lust
- This COS top from the men's section ticks all the right...**: Carlo Bevelander Low Top
- stan smith outfits ~ Buscar con Google**: Denys Finch-Hatton Sneakers
- Glorious Ladies**: Glorious Ladies

Building object-to-scene visual search leveraging large-scale product detection



Improving semantic object detectors





Multi-modal search

with no intent information from user, we must infer a **use-case**



Multi-modal search

with no intent information from user, we must infer a **use-case**

table

Decorations Settings Coffee Design

See more

Buyable Pins
Buy on Pinterest

\$188
Reaching Branches Table Lamp
GoGetGlam
 GLAM NYC -...
Home Decor -...

\$550
Granada Side Table
CROFT HOUSE
 Croft House
TABLES | C...

\$350
Industrial Table Leg
Bold MFG & Co.
 Bill You
Details

Conclusions

- Deep learning = powerful **featurizer**
- Object detectors = product co-occurrence signals
- Search products moving towards **multi-modal search**

