# Self-Supervised Fine-Grained Food Recognition
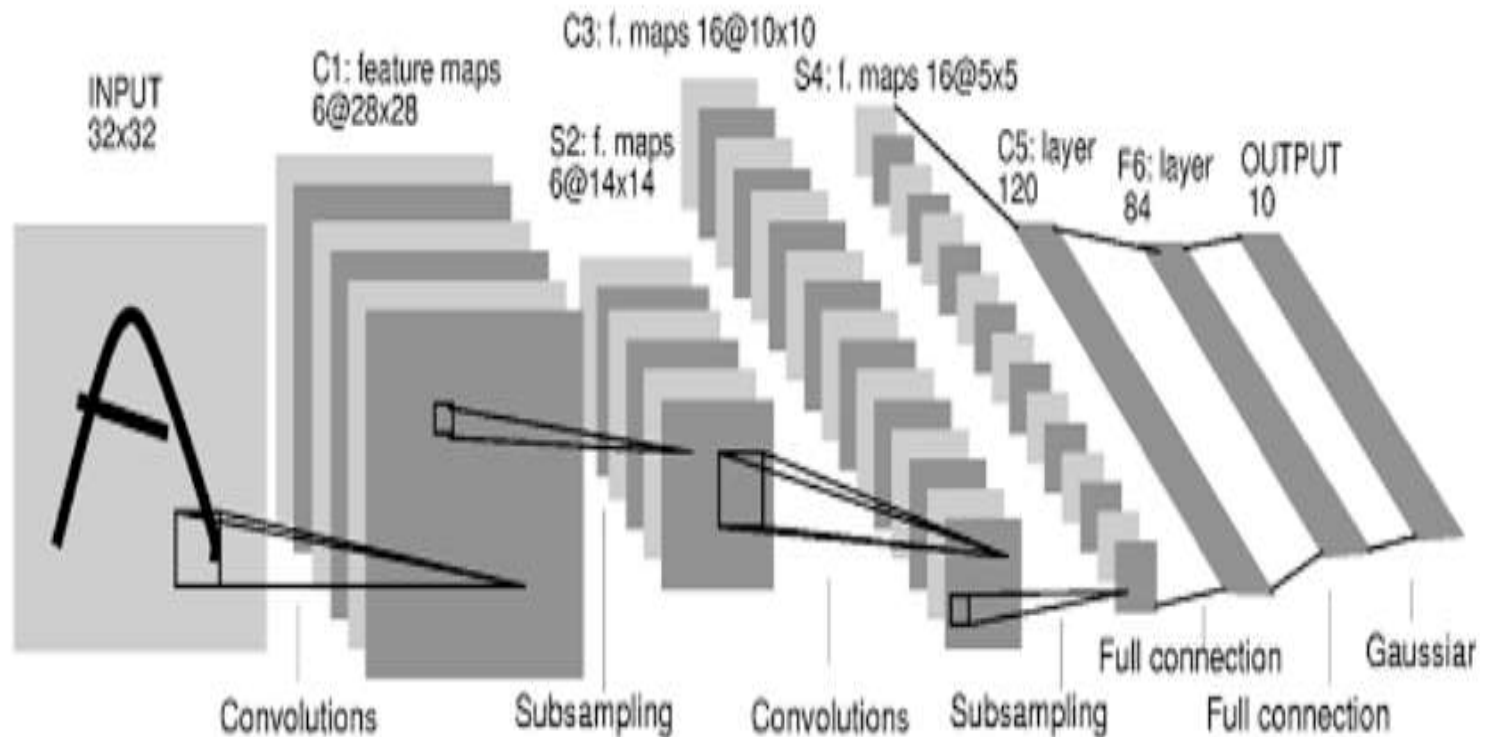


Petia Radeva
Universitat de Barcelona &
Computer Vision Center, Spain



UNIVERSITAT DE BARCELONA



CVC
Centre de Visió per Computador

ICPRAM 2023
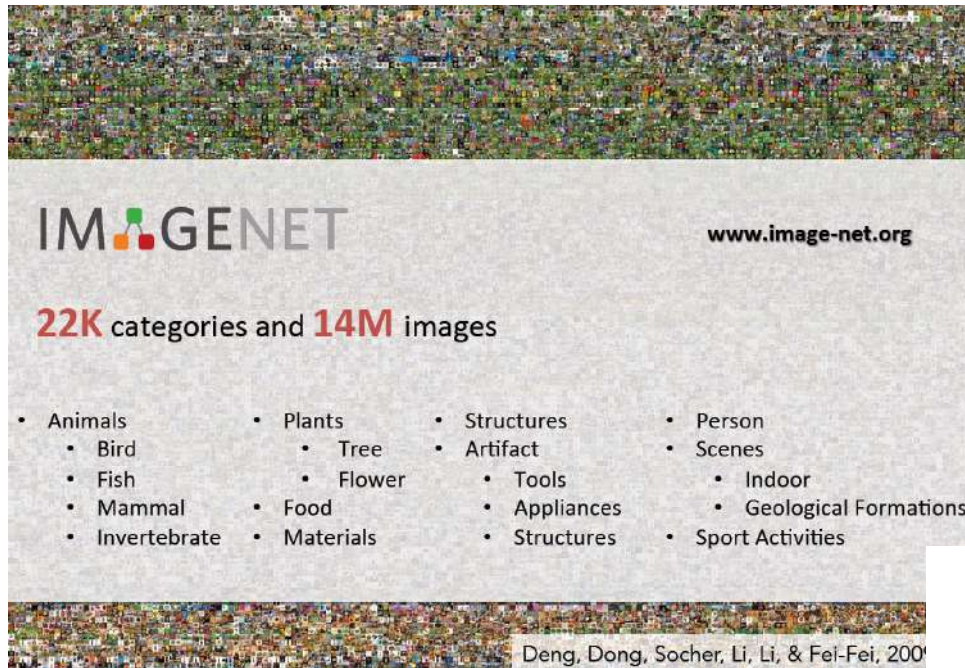12th International Conference on Pattern Recognition Applications and Methods

1998

LeCun et al.

LeCun, Yann; Léon Bottou; Yoshua Bengio; Patrick Haffner (1998). "Gradient-based learning applied to document recognition". *Proceedings of the IEEE* **86** (11): 2278–2324

# Imagenet



**IMAGENET**

www.image-net.org

**22K** categories and **14M** images

- Animals
  - Bird
  - Fish
  - Mammal
  - Invertebrate
- Plants
  - Tree
  - Flower
- Food
- Materials
- Structures
- Artifact
  - Tools
  - Appliances
  - Structures
- Person
- Scenes
  - Indoor
  - Geological Formations
- Sport Activities

Deng, Dong, Socher, Li, Li, & Fei-Fei, 2009



TED Ideas worth spreading          WATCH   DISCOVER   ATT

Fei-Fei Li:

**How we're teaching computers to understand pictures**
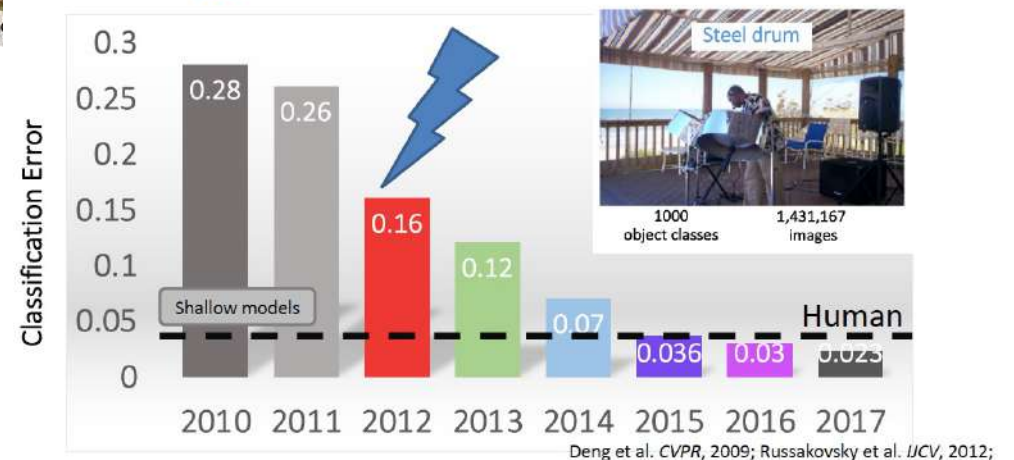
TED2015 · 17:58 · Filmed Mar 2015
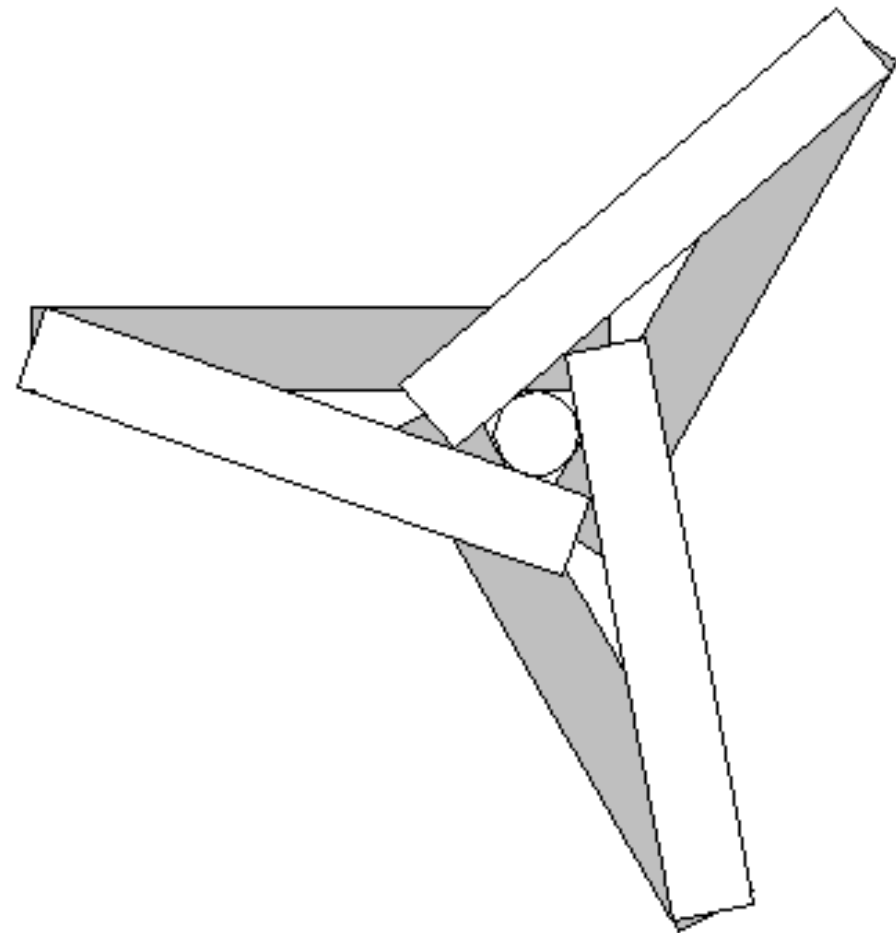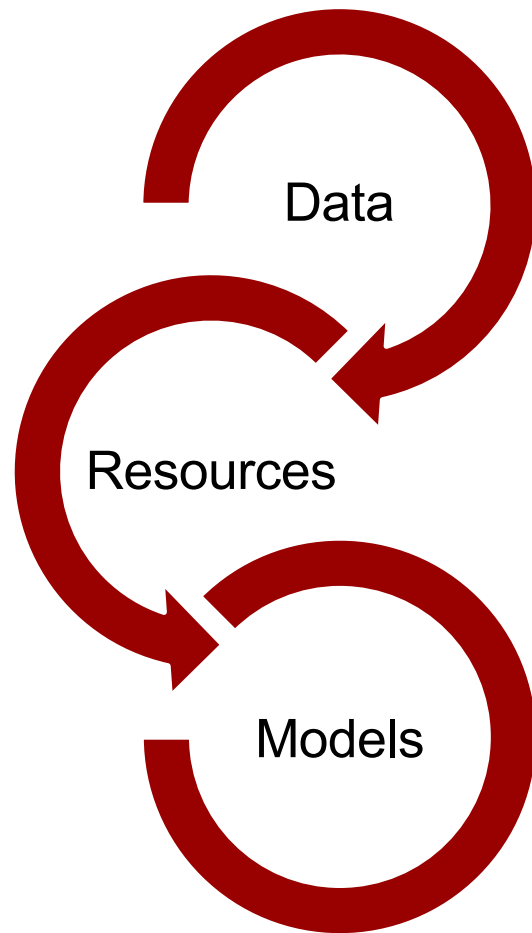
26 subtitle languages

View interactive transcript

**1,607,730** Total views

**IMAGENET Classification Task**



Steel drum

1000 object classes     1,431,167 images

Classification Error

0.3
0.25    0.28   0.26
0.2
0.15              0.16
0.1                      0.12
0.05   Shallow models          0.07    Human
0                              0.036  0.03  0.023

2010 2011 2012 2013 2014 2015 2016 2017

Deng et al. *CVPR*, 2009; Russakovsky et al. *IJCV*, 2012;

21:22

Data

Resources

Models

# The Importance of GPUs

- Nvidia Tensor Cores - 2017

- Google Tensor Processing Unit (TPU) - 2016

- Intel - Nervana Neural Processor - 2017

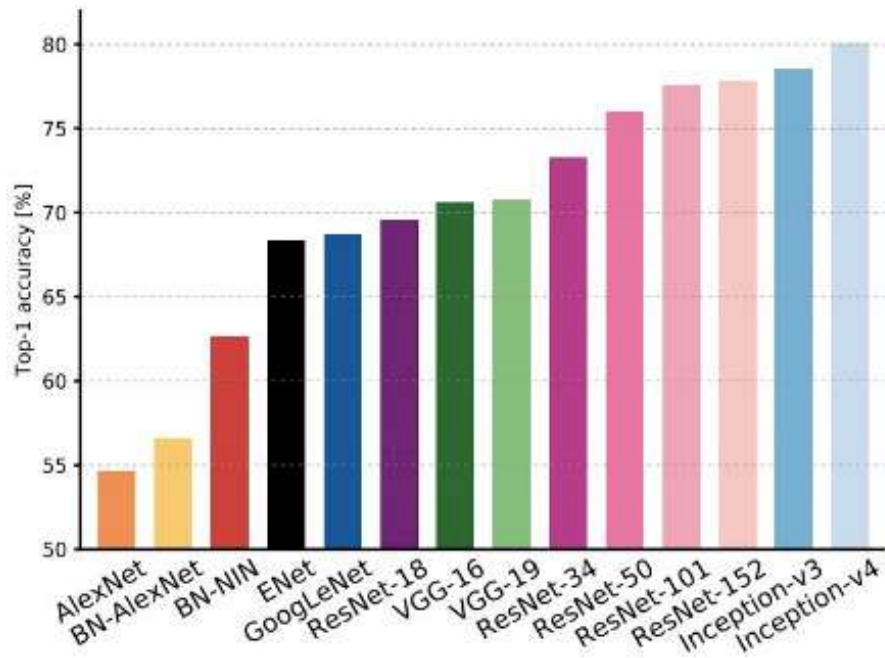- GPUs in Cloud Computing (Google, 2017)



$$D = \begin{pmatrix} A_{0,0} & A_{0,1} & A_{0,2} & A_{0,3} \\ A_{1,0} & A_{1,1} & A_{1,2} & A_{1,3} \\ A_{2,0} & A_{2,1} & A_{2,2} & A_{2,3} \\ A_{3,0} & A_{3,1} & A_{3,2} & A_{3,3} \end{pmatrix} \begin{pmatrix} B_{0,0} & B_{0,1} & B_{0,2} & B_{0,3} \\ B_{1,0} & B_{1,1} & B_{1,2} & B_{1,3} \\ B_{2,0} & B_{2,1} & B_{2,2} & B_{2,3} \\ B_{3,0} & B_{3,1} & B_{3,2} & B_{3,3} \end{pmatrix} + \begin{pmatrix} C_{0,0} & C_{0,1} & C_{0,2} & C_{0,3} \\ C_{1,0} & C_{1,1} & C_{1,2} & C_{1,3} \\ C_{2,0} & C_{2,1} & C_{2,2} & C_{2,3} \\ C_{3,0} & C_{3,1} & C_{3,2} & C_{3,3} \end{pmatrix}$$
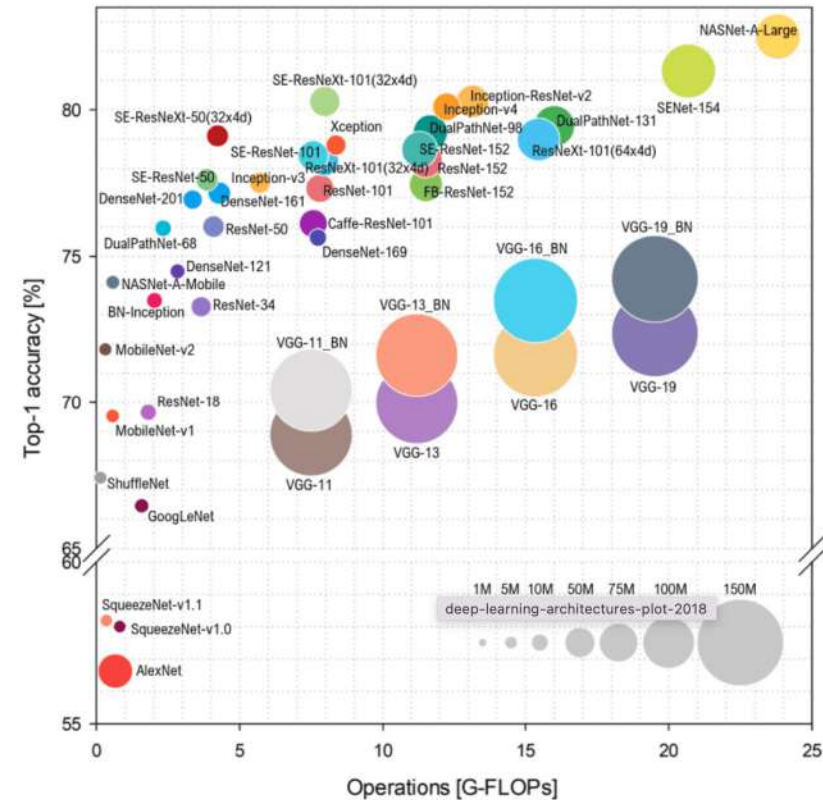
FP16 or FP32      FP16      FP16      FP16 or FP32

GPU cores is based on matrix multiplication

21:22

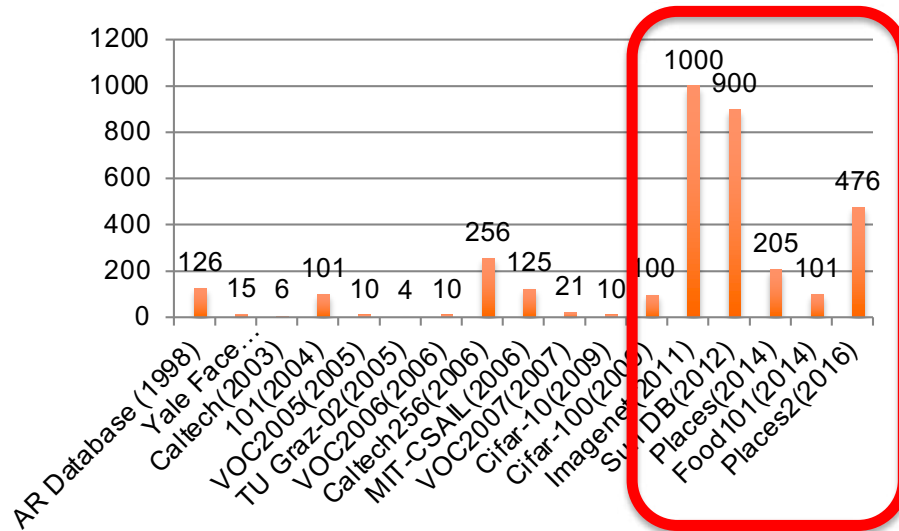https://www.doc.ic.ac.uk/~jce317/history-machine-learning.html#top

- **Millions of parameters!!!**

The process of training a CNN consists of training all hyperparameters: convolutional matrices and weights of the fully connected layers.
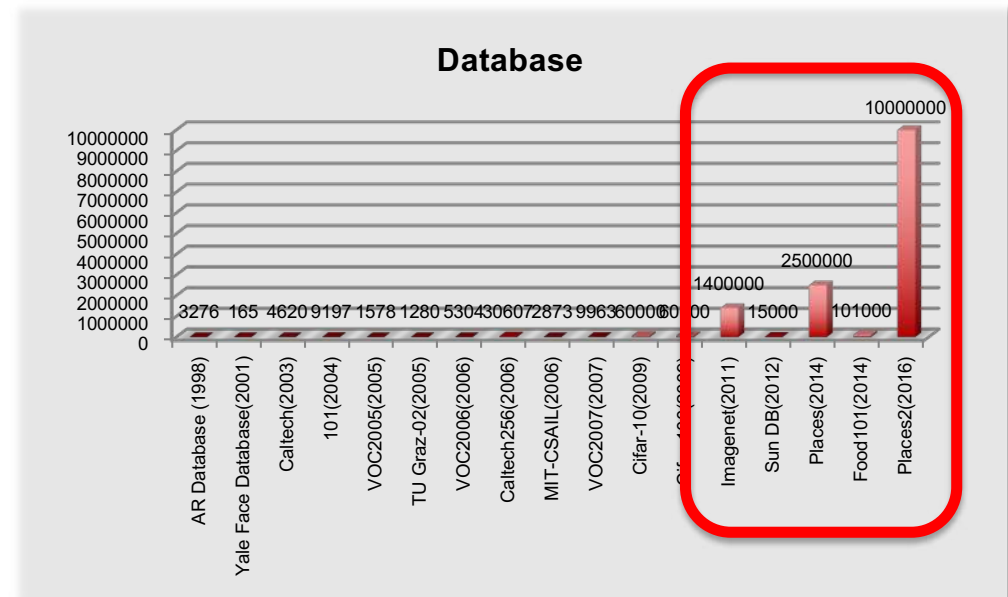
Number of objects/Database



**ImageNet & Deep learning**

Number of images/Database



21:22

# What are the most popular datasets today?

| Dataset | Papers | Benchmarks | Images (K) | Classes | Sizes |
|---|---|---|---|---|---|
| Cifar-10 | 10581 | 66 | 60 | 10 | 32x32 |
| **ImageNet** | 10046 | 97 | **1400** | **1000** | variable |
| COCO | 7160 | 78 | 123 | 80 | |
| MNIST | 5911 | 49 | 60 | 10 | 28x28 |
| Cifar-100 | 5322 | 42 | 60 | 100 | 32x32 |
| Cityskapes | 2562 | 37 | 25 | 8 | |
| SVHN | 2474 | 11 | 60 | 10 | 32x32 |
| Kitti | 2453 | 120 | 0,5 | 11 | |
| CelebA | 2408 | 20 | 202 | 10177 | 178x218 |
| Fashion-MNIST | 2150 | 17 | 70 | 10 | 28x28 |
| CUB-00-2011 | 2408 | 37 | 12 | 200 | |
| **Places** | 760 | 4 | **2500** | 205 | |
| Tiny ImageNet | 516 | 7 | 31 | 200 | |
| **Places205** | 468 | 1 | **2500** | 205 | |
| Caltech-101 | 393 | 6 | 5 | 101 | 300x200 |
| Stanford Cars | 392 | 8 | 16 | 196 | 360x240 |
| Caltech-256 | 345 | 4 | 30 | 257 | |

# Large Scale Food Recognition Dataset

## Large Scale Visual Food Recognition

Publisher: IEEE    Cite This    PDF

Weiqing Min ; Zhiling Wang ; Yuxin Liu ; Mengjiang Luo ; Liping Kang ; Xiaoming Wei ; Xia...    All Authors
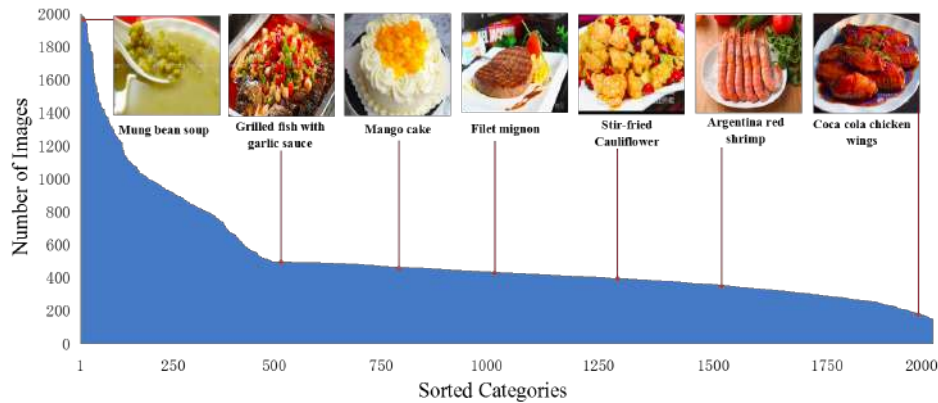
### Abstract

Abstract:
Food recognition plays an important role in food choice and intake, which is essential to the health and well-being of humans. It is thus of importance to the computer vision community, and can further support many food-oriented vision and multimodal tasks, e.g., food detection and segmentation, cross-modal recipe retrieval and generation.

Authors

Keywords

Metrics

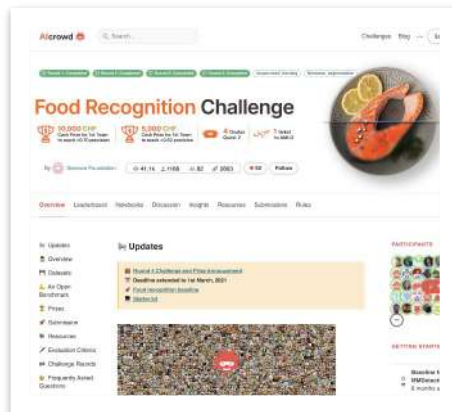| Meat | Cheese back ribs | Tomahaw | Fried pork in scoop | Sheep roll |
|---|---|---|---|---|
| Vegetables | Eggplant salad | Fruit salad | Shredded cucumber | Fried eggplant |
| Bread | Tuna pizza | Beef burger | Seafood pancake | Coconut bread |
| Snack | Egg tart | Roti prata | Strawberry smoothie | Takoyaki |
| Fried food | Tonkatsu | Fried chicken | Fried cuttlefish balls | Fried tofu |
| Seafood | Tempura | Spicy crab | Geoduck sashimi | Cod fish steak |
| Cereal products | Egg fried rice | Salmon sushi | Pan-fried pork bun | Instant noodles |

# Food recognition popularity



## Number of Food recognition papers





iFood 2011 fine-grained (prepared) food categories with 135733
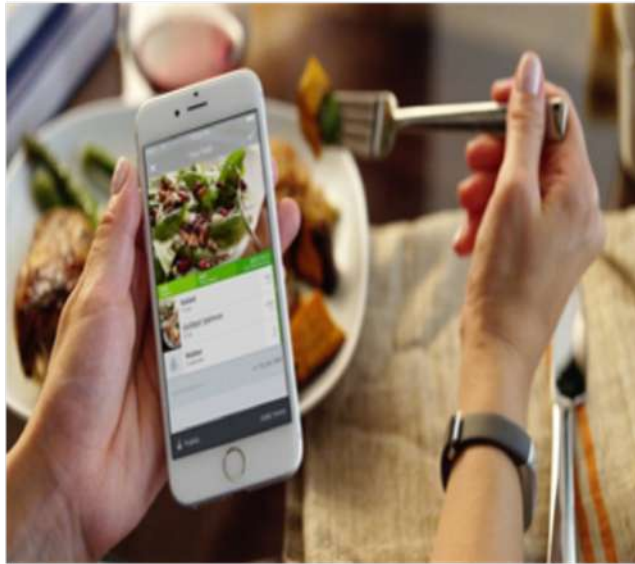


AICrowd: 26000 annotated segmented images



LargeFineFoodAI: 1,000 fine-grained food categories and over 50,000 images.

Food image analysis

# Why food recognition?



**"Camera eats first"**

180M #food
90/minute





54% take picture
39% post it

21:22

# Why is the food recognition a challenge?

# Food Analysis Problems

Ingredients
- Intra-class variability

- Inter-class similarity



*Intra-class variability example: Apple. Image source: Recipes5k*



*Inter-class similarity example: Tomato sauce and Curry sauce. Image source: Recipes5k*

# Decreasing in Precision

21:22

The food recognition is a Fine-grained recognition problem

# Challenges of Food image analysis

Food256: 25.600 images (100 images/class) Classes: 256



Food101 – 101.000 images
(1000 images/class)
Classes: 101

FoodX-251
Classes: 251
140K images

Food1K
Classes: 1000
370K images

**Food DB**

150.000 images
231 categories

**ImageNet**

1.400.000 images
1000 categories

**Future Food DB**

????? images
200.000 categories

**Current SoA on Food recognition**
- 79% on UECFOOD
- 44% on ChinaFood1000

21:22

# SSL: Benefits & Uses in 2023



https://research.aimultiple.com/self-supervised-learning/

# Self-Taught AI Shows Similarities to How the Brain Works

SSL allows a neural network to **figure out for itself what matters**.

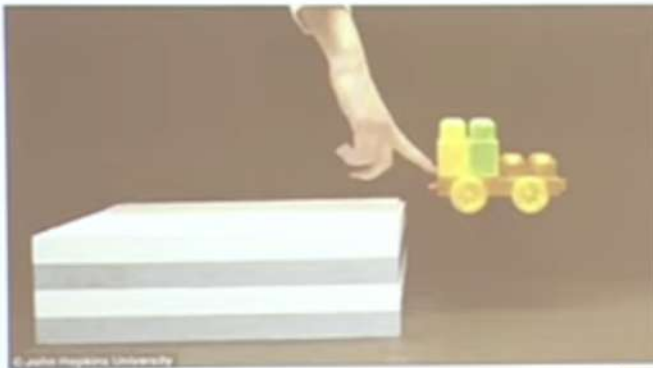Explore neural networks trained **with little or no human-labelled data.**

Computational **models of the mammalian visual and auditory systems** built using self-supervised learning models have shown a **closer correspondence to brain** function than their supervised-learning counterparts.



Alexei Efros, University of California, Berkeley, "Most modern AI systems are too reliant on human-created labels. They don't really learn the material".

https://www.quantamagazine.org/self-taught-ai-shows-similarities-to-how-the-brain-works-20220811/

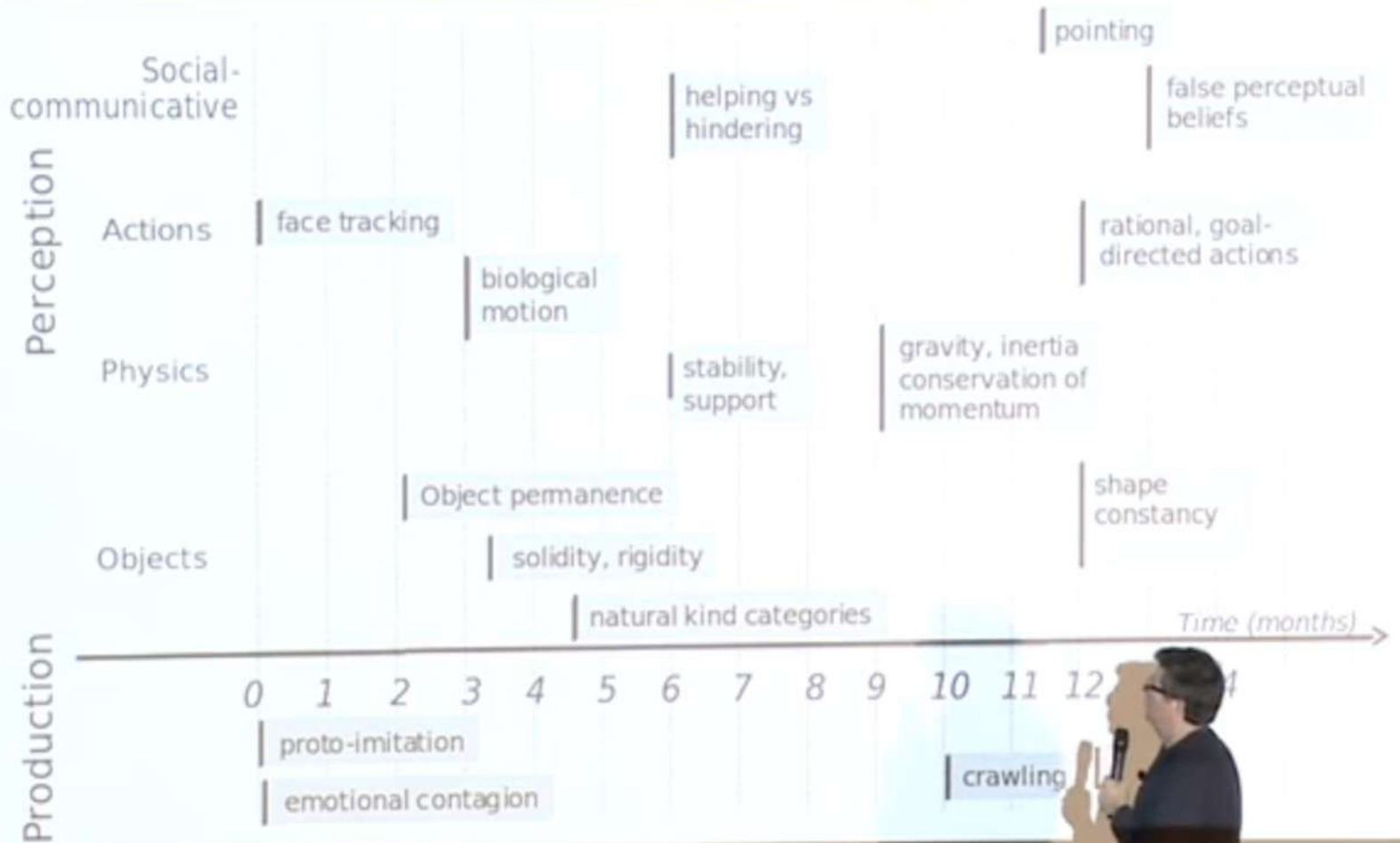# Babies learn how the world works by observation



▶ Largely by observation, with remarkably little interaction.

Photos courtesy of
Emmanuel Dupoux

21:22

# Early conceptual acquisition in infants (from Emmanuel Dupoux)

**Understand brain through NNs:**

- the brain is full of feedback connections, while current models have few such connections, if any.

**An obvious next step:** use SSL to train highly recurrent networks and see how the activity in NNs compares to real brain activity.

**Crucial step:** match the activity of NNs in SSL models to the activity of individual biological neurons.



"No doubt that 90% of what the brain does is self-supervised learning," Blake Richards, a computational neuroscientist at McGill University and Mila, the Quebec Artificial Intelligence Institute.

**Hypothesis**: the visual systems of humans and other primates are the best studied of all animal sensory systems,
- but neuroscientists have struggled to explain why they include two separate pathways:
  - the ventral visual stream, which is responsible for recognizing objects and faces, and
  - the dorsal visual stream, which processes movement (the "what" and "where" pathways, respectively).

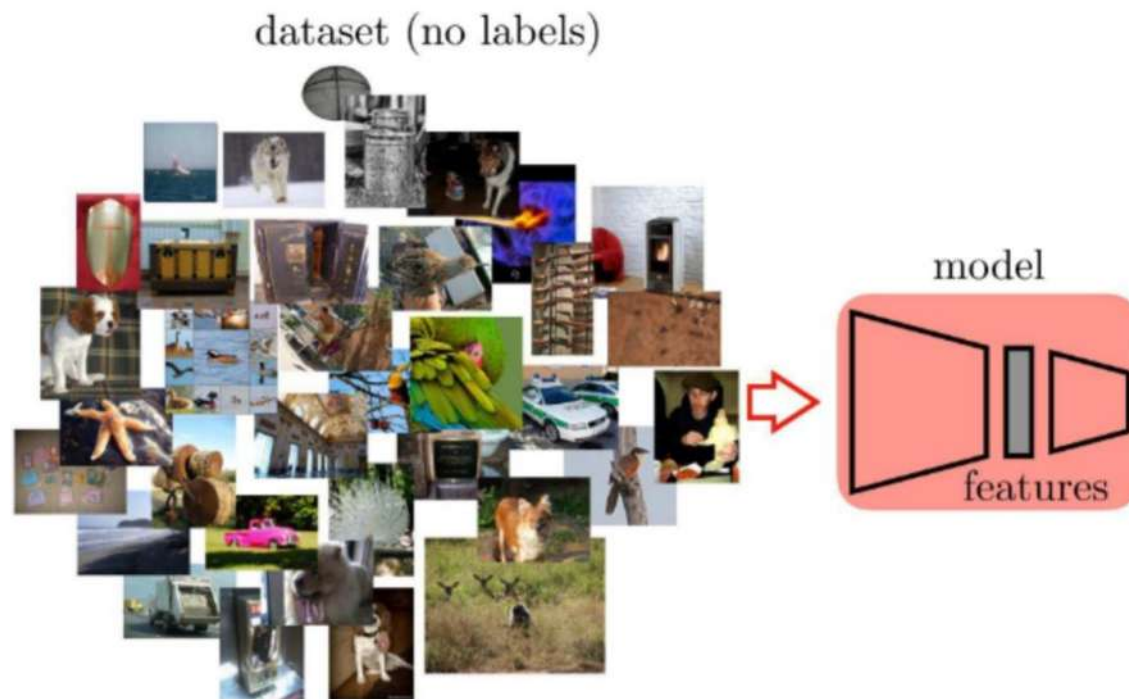# Self-supervised vs supervised learning



JAISWAL, Ashish, et al. A survey on contrastive self-supervised learning.Technologies, 2020, 9.1: 2.

dataset (no labels)

Contrastive loss

$$\mathcal{L}_i^{\text{InfoNCE}} = -\log \frac{\exp\left(z_i \cdot z_i^+/\tau\right)}{\exp\left(z_i \cdot z_i^+/\tau\right) + \sum_{z^- \in \mathcal{N}_i} \exp\left(z_i \cdot z^-/\tau\right)}$$

**MoCo**

**SimCLR**

**SwAV**

**BYOL**

**SimSiam**

**Barlow Twins**

21:22

# Momentum Contrasting (MOCO)

Given an image $x_i$, MoCo learns a query encoder $q = f_q(x_i)$ able to **differentiate $q_i = f_q(x_i)$ from the other images**.

**Positive pairs**: 2 representations of the same image without augmentation.

> An asynchronously updated momentum encoder $f_k(.)$ is used to generate the positive counterpart $k^+ = f_k(x_i)$.

**Negative samples:** MoCo derives from a memory bank, storing previously encoded representations.

The model optimizes the following objective function:

$$L_i^{MoCo} = -log\left(\frac{exp(q_i \cdot k_i^+/\tau)}{\sum_{k=1}^{K} exp(q_i \cdot k_k^-/\tau)}\right)$$

where K is the number of negative samples in the queue.



Kaiming He et al. "Momentum Contrast for Unsupervised Visual Representation Learning", CVPR 2020, pp. 9729–9738.

EMA denotes exponential moving average updates.



$$\xi \leftarrow \tau\xi + (1 - \tau)\theta$$

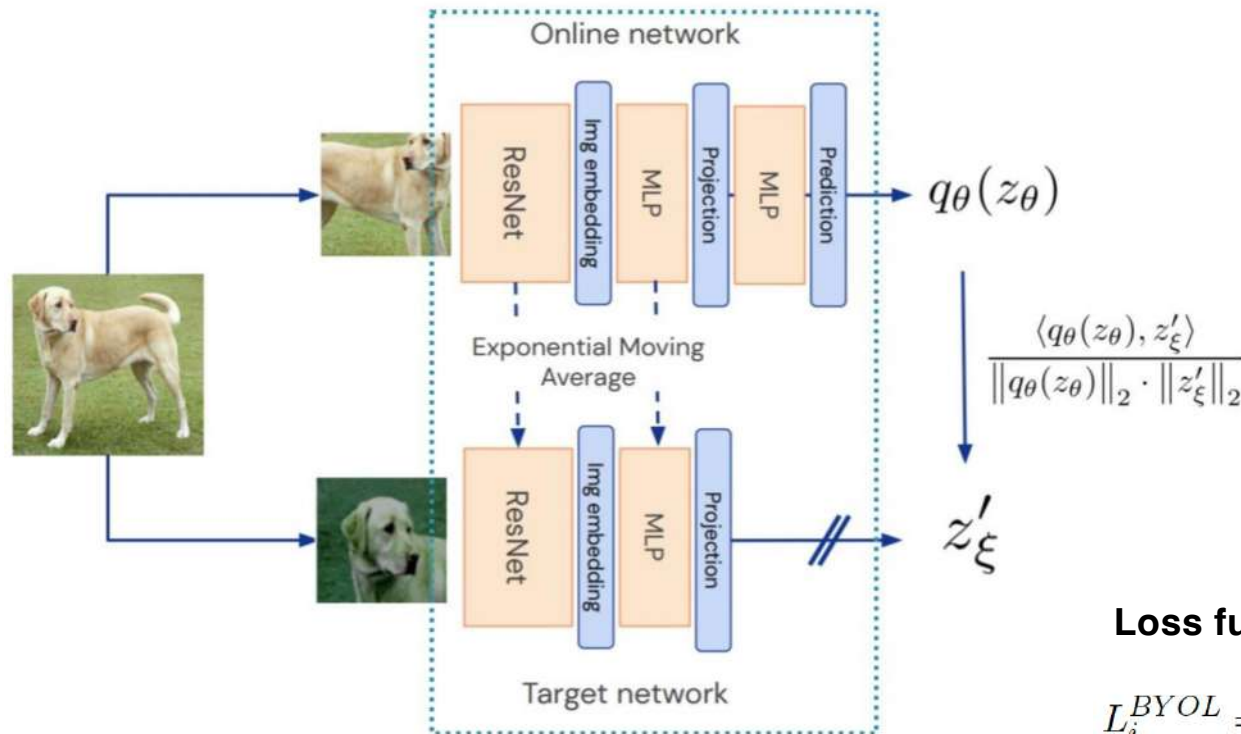The EMA-teacher framework with standard Batch Normalization.

$$\mathcal{L}=\frac{1}{2}\mathcal{D}(p_1, \text{stopgrad}(z_2))+\frac{1}{2}\mathcal{D}(p_2, \text{stopgrad}(z_1))$$



**NO:**
(i) negative sample pairs,
(ii) large batches,
(iii) momentum encoders.

X. Chen and K. He. "Exploring Simple Siamese Representation Learning". CVPR, 2021, pp. 15750–15758

**Avoid negative sampling!**

**Introduces predictor:**

It tries to predict different views (regression targets) of the same image directly in the representation space using a predictor.

**Loss function:**

$$L_i^{BYOL} = \left\| p_i^1 - z_i^2 \right\|_2^2 = 2 - 2 \cdot \frac{\langle p_i^1, z_i^2 \rangle}{\left\| p_i^1 \right\|_2 \cdot \left\| z_i^2 \right\|_2}.$$

J.-B. Grill et al. "Bootstrap Your Own Latent - A New Approach to Self-Supervised Learning". NIPS. Vol. 33. 2020, pp. 21271–21284

# SimCLR by the Google AI team

**Introduces projectors:** a learnable nonlinear transformation between the representation and the contrastive loss

**Positive sampling**: Given a batch of N samples, the pretext task P generates two augmented views $x^a_i$ and $x^+_i$ for each sample $x_i$ of the batch.

**Negative sampling**: the rest of the images $x^-_i$ on the same batch to form the negative pairs $(x^a_i, x^-_i)$.

Batch sizes of 8196 are used.

**Loss function:**

$$L_i^{SimCLR} = -log\left(\frac{exp(z_i^a \cdot z_i^+/\tau)}{\sum_{k=1}^{N} exp(z_i^a \cdot z_k^-/\tau)}\right)$$

Ting Chen et al. "A Simple Framework for Contrastive Learning of Visual Representations". 37th ICMLL 2020, pp. 1597–1607.

https://analyticsindiamag.com/what-is-contrastive-self-supervised-learning/

$$\left(z_i, z_i^+\right)$$

Maximize agreement

$z_i \longleftrightarrow z_j$

Proj. head $\longrightarrow g(\cdot)$   $g(\cdot) \longleftarrow$ Proj. head

$\boldsymbol{h}_i$   $\longleftarrow$ Representation $\longrightarrow$   $\boldsymbol{h}_j$

ResNet50 $\longrightarrow f(\cdot)$   $f(\cdot) \longleftarrow$ ResNet50

$\tilde{\boldsymbol{x}}_i$   $\tilde{\boldsymbol{x}}_j$

$t \sim \mathcal{T}$   $\boldsymbol{x}$   $t' \sim \mathcal{T}$

SimCLR

manifold of all samples

nearest neighbor of view 1 in the support set

$$\mathcal{L}_i^{\text{NNCLR}} = -\log \frac{\exp\left(\text{NN}(z_i, Q) \cdot z_i^+ / \tau\right)}{\sum_{k=1}^{n} \exp\left(\text{NN}(z_i, Q) \cdot z_k^+ / \tau\right)}$$

Mean shift for Self-supervised learning (MSF)



$$L_i^{MSF} = \frac{1}{K} \sum_{k=1}^{K} dist(nn_i^1 k, p_i^2)$$

- Improve the retrieval without efficiency loss
- Efficiency decrease

S. A. Koohpayegani, A. Tejankar, and H. Pirsiavash. "Mean Shift for Self-Supervised Learning", CVPR 2021, pp. 10326–10335.

21:22

J. Zbontar et al. "Barlow Twins: Self-Supervised Learning via Redundancy Reduction". 38th ICML. 2021, pp. 12310–12320

21:22

## Barlow Twins architecture

## Barlow Twins architecture

## Barlow Twins architecture

## Barlow Twins architecture

C

|       | $f_1$ | $f_2$ | $f_3$ | $f_4$ | $f_5$ |
|-------|-------|-------|-------|-------|-------|
| $f_1$ |       |       |       |       |       |
| $f_2$ |       |       |       |       |       |
| $f_3$ |       |       |       |       |       |
| $f_4$ |       |       |       |       |       |
| $f_5$ |       |       |       |       |       |

**Barlow Twins' Loss Function**

$$\mathcal{L}_{\mathcal{BT}} = \sum_{i=1}^{D} (1 - c_{ii})^2 + \lambda \sum_{i=1}^{D} \sum_{j=1,j\neq i}^{D} c_{ij}^2$$

Invariance term      Redundancy reduction term

21:22

Remember: in the SimSiam architecture

Remember: in the SimSiam architecture

Remember: in the SimSiam architecture
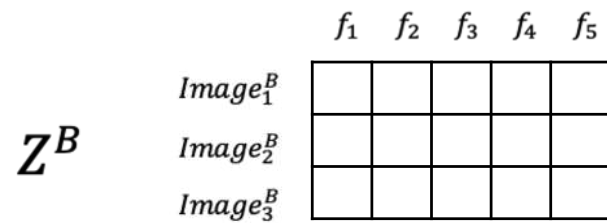
Remember: in the SimSiam architecture

## SimSiam architecture



$$\mathcal{L}_{\text{SimSiam}} = -\frac{1}{2N}\sum_{i=1}^{N}\left(\left(Sim_i^{(1)}\right)^2 + \left(Sim_i^{(2)}\right)^2\right)$$

21:22

$$C^{(1)}$$

$$\begin{array}{c} \\ I_1 \\ I_2 \\ I_3 \end{array}$$ (with $I_1\ I_2\ I_3$ column headers)

$$C^{(2)}$$

$$\begin{array}{c} \\ I_1 \\ I_2 \\ I_3 \end{array}$$ (with $I_1\ I_2\ I_3$ column headers)

**OptSSL's Loss Function**

$$\mathcal{L}_{Opt\text{-}SSL} = \mathcal{L}_{i\text{-}diag} + \lambda_1 \cdot \mathcal{L}_{i\text{-}off\text{-}diag} + \mathcal{L}_{f\text{-}diag} + \lambda_2 \cdot \mathcal{L}_{f\text{-}off\text{-}diag}$$

where

$$\mathcal{L}_{diag} = \sqrt{\frac{1}{2N}\left(\sum_{i=1}^{N}\left(1 - c_{ii}^{(1)}\right)^2 + \sum_{i=1}^{N}\left(1 - c_{ii}^{(2)}\right)^2\right)}$$

$$\mathcal{L}_{off\text{-}diag} = \sqrt{\frac{1}{2N(N-1)}\left(\sum_{i=1}^{N}\sum_{j=1,j\neq i}^{N}\left(c_{ij}^{(1)}\right)^2 + \sum_{i=1}^{N}\sum_{j=1,j\neq i}^{N}\left(c_{ij}^{(2)}\right)^2\right)}$$

**Applied both**: to images in the batch and the features!

**OptSSL Architecture**

Identity matrix

Similarity

Cross-correlation matrix

grad

Stop-grad

Batch dimension

$P^A$

Predictor h

$Z^A$      $Z^B$

Encoder $f$      Encoder $f$

$X^A$   $t \sim \mathcal{T}$    $t' \sim \mathcal{T}$   $X^B$

$\mathcal{P}$

$X$

**Encoder $f$**

224 x 224 x 3

Backbone (ResNet - 50)

2048

$fc$ layer 1

Batch norm.

ReLU

2048

$fc$ layer 2

Batch norm.

ReLU

2048

$fc$ layer 3

Batch norm.

2048

**Predictor h**

2048

$fc$ layer 1

Batch norm.

ReLU

512

$fc$ layer 2

2048

Nil Ballús, Bhalaji Nagarajan, Petia Radeva: Opt-SSL: An Enhanced Self-Supervised Framework for Food Recognition. IbPRIA 2022: 655-666

21:22

# Reducing Redundancy: Feature Contrast



$$Z^1 = \begin{bmatrix} z^1_{1\,1} & z^1_{1\,2} & \cdots & z^1_{1\,256} \\ z^1_{2\,1} & z^1_{2\,2} & \cdots & z^1_{2\,256} \\ \vdots & \vdots & \ddots & \vdots \\ z^1_{N\,1} & z^1_{N\,2} & \cdots & z^1_{N\,256} \end{bmatrix}$$

$$Z^2 = \begin{bmatrix} z^2_{1\,1} & z^2_{1\,2} & \cdots & z^2_{1\,256} \\ z^2_{2\,1} & z^2_{2\,2} & \cdots & z^2_{2\,256} \\ \vdots & \vdots & \ddots & \vdots \\ z^2_{N\,1} & z^2_{N\,2} & \cdots & z^2_{N\,256} \end{bmatrix}$$

Similar

Momentum Encoder

$g_\xi$

EMA

Online Encoder

$g_\theta$

11
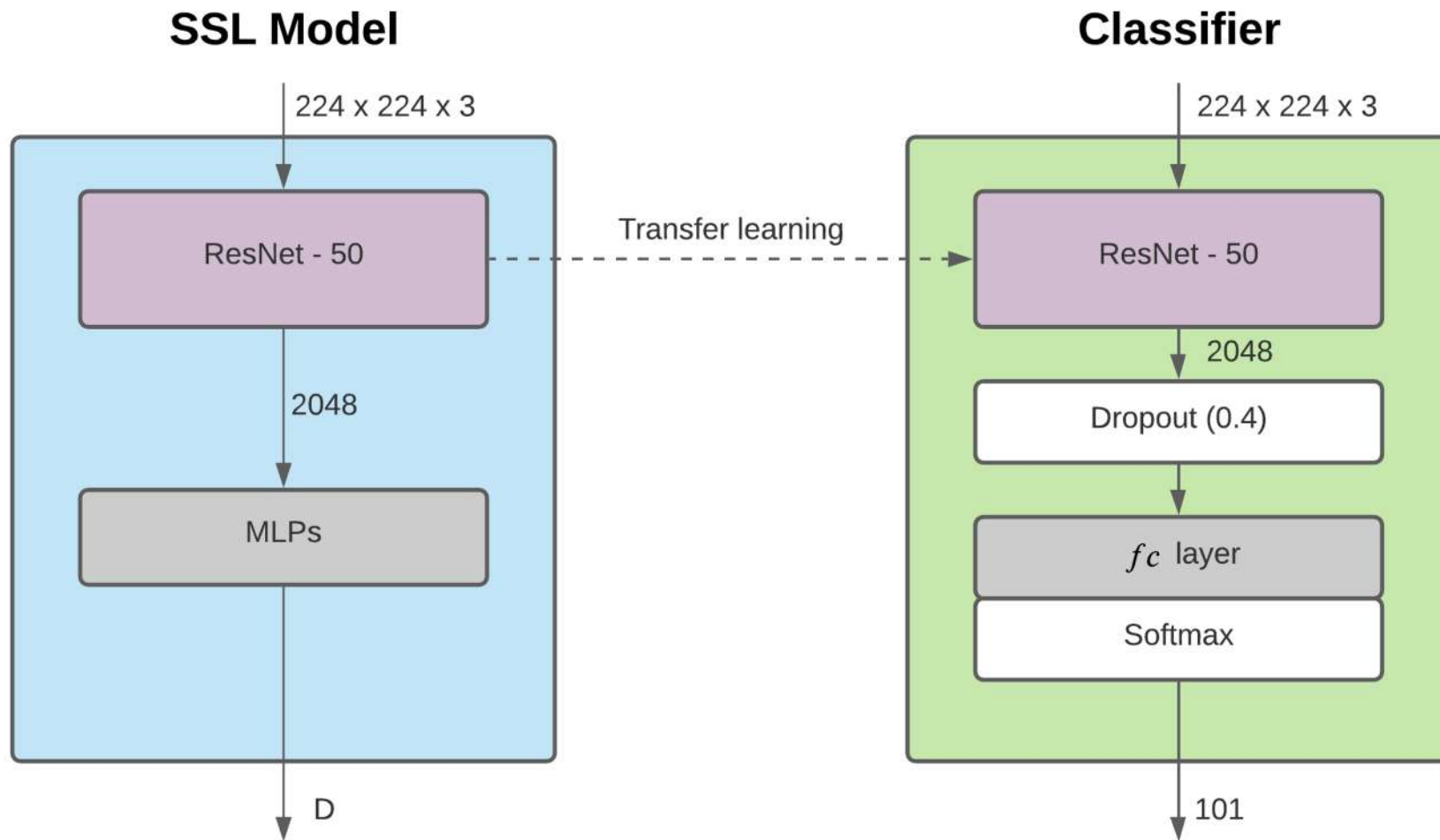
# Validation

## Dataset



Food-101

## Evaluation Metrics

### SSL Model

- Top 1 and top 5 accuracy using a k-NN classifier

### Classifier

- Overall accuracy
- Variance
- Entropy
- Mutual Information

21:22

Performance of different SSL Models



| SSL Model | λ | K-NN Acc. | |
|-----------|-----|------------------|------------------|
| | | Top 1 acc. (%) | Top 5 acc. (%) |
| Barlow Twins | 0,0051 | 49,89 | 79,07 |
| SimSiam | - | 54,33 | 81,62 |
| OptSSL | 0,5 | **63,52** | **87,48** |

21:22

## SSL Frameworks: Results

| Framework | Main Characteristics |
|---|---|
| SimCLR | Starting point, use of Negative samples, high batch size, contrastive loss |
| BYOL | Randomly initialized network instead of Negative samples, MSE loss, lower reliance on batch size, introduces the predictor |
| SimSiam | SimCLR without Negative samples, BYOL without the target network, use of stop-gradient, siamese structure |
| Barlow Twins | Innovative loss function that uses cross-correlation matrixes, use of batch normalization, benefits from high dimensional representations |
| NNCLR | Introduces NN algorithm to provide more richness, modified version of SimCLR |

| Method | Top-1 | Top-5 |
|---|---|---|
| PIRL | 63.6 | - |
| CPC v2 | 63.8 | 85.3 |
| PCL | 65.9 | - |
| CMC | 66.2 | 87.0 |
| MoCo v2 | 71.1 | - |
| SimSiam | 71.3 | - |
| SimCLR v2 | 71.7 | - |
| SwAV | 71.8 | - |
| InfoMin Aug. | 73.0 | 91.1 |
| BYOL | 74.3 | 91.6 |
| NNCLR (ours) | **75.4** | **92.3** |
| BARLOW TWINS (ours) | 73.2 | 91.0 |

ImageNet linear classification results from Debidatta Dwibedi et al. (2021)

21:22

# Quantitative Results: CIFAR

CIFAR-10

| Method | Backbone | Epochs | Acc@1 (Online) | Acc@5 (Online) |
|---|---|---|---|---|
| BYOL | ResNet18 | 1000 | 92.58 | 99.79 |
| DeepCluster V2 | ResNet18 | 1000 | 88.85 | 99.58 |
| DINO | ResNet18 | 1000 | 89.52 | 99.71 |
| MoCo V2+ | ResNet18 | 1000 | 92.94 | 99.79 |
| MoCo V3 | ResNet18 | 1000 | 93.10 | 99.80 |
| ReSSL | ResNet18 | 1000 | 90.63 | 99.62 |
| SimCLR | ResNet18 | 1000 | 90.74 | 99.75 |
| Simsiam | ResNet18 | 1000 | 90.51 | 99.72 |
| SwAV | ResNet18 | 1000 | 89.17 | 99.68 |
| VIbCReg | ResNet18 | 1000 | 91.18 | 99.74 |
| VICReg | ResNet18 | 1000 | 92.07 | 99.74 |
| W-MSE | ResNet18 | 1000 | 88.67 | 99.68 |
| Barlow Twins | ResNet18 | 1000 | 92.10 | 99.73 |
| NNCLR | ResNet18 | 1000 | 91.88 | 99.78 |
| **Musketeer (Ours)** | ResNet18 | 1000 | **93.24** | **99.88** |

CIFAR-100

| Method | Backbone | Epochs | Acc@1 (Online) | Acc@5 (Online) | k-NN Acc@1 (Online) |
|---|---|---|---|---|---|
| BYOL | ResNet18 | 1000 | 70.46 | 91.96 | - |
| DeepCluster V2 | ResNet18 | 1000 | 63.61 | 88.09 | - |
| DINO | ResNet18 | 1000 | 66.76 | 90.34 | - |
| MoCo V2+ | ResNet18 | 1000 | 69.89 | 91.65 | - |
| MoCo V3 | ResNet18 | 1000 | 68.83 | 90.57 | - |
| ReSSL | ResNet18 | 1000 | 65.92 | 89.73 | - |
| SimCLR | ResNet18 | 1000 | 65.78 | 89.04 | - |
| Simsiam | ResNet18 | 1000 | 66.04 | 89.62 | - |
| SwAV | ResNet18 | 1000 | 64.88 | 88.78 | - |
| VIbCReg | ResNet18 | 1000 | 67.37 | 90.07 | - |
| VICReg | ResNet18 | 1000 | 68.54 | 90.83 | - |
| W-MSE | ResNet18 | 1000 | 61.33 | 87.26 | - |
| NNCLR | ResNet18 | 1000 | 69.62 | 91.52 | - |
| NNCLR* | ResNet18 | 1000 | 69.17 | 91.70 | 62.16 |
| Barlow Twins | ResNet18 | 1000 | 70.90 | 91.91 | - |
| Barlow Twins* | ResNet18 | 1000 | 71.21 | 92.46 | 63.11 |
| MSF* | ResNet18 | 1000 | 67.84 | 91.64 | 63.36 |
| **Musketeer (Ours)** | ResNet18 | 1000 | **72.17** | **93.35** | **64.84** |

# Quantitative Results: ImageNet-100

| Method | Backbone | Epochs | Acc@1 (online) | Acc@5 (online) |
|---|---|---|---|---|
| BYOL ++ | ResNet18 | 400 | 80.16 | 95.02 |
| DeepCluster V2 | ResNet18 | 400 | 75.36 | 93.22 |
| DINO | ResNet18 | 400 | 74.84 | 92.92 |
| MoCo V2+ ++ | ResNet18 | 400 | 78.20 | 95.50 |
| MoCo V3 ++ | ResNet18 | 400 | 80.36 | 95.18 |
| ReSSL | ResNet18 | 400 | 76.92 | 94.20 |
| SimCLR ++ | ResNet18 | 400 | 77.64 | 94.06 |
| Simsiam | ResNet18 | 400 | 74.54 | 93.16 |
| SwAV | ResNet18 | 400 | 74.04 | 92.70 |
| VIbCReg | ResNet18 | 400 | 79.86 | 94.98 |
| VICReg ++ | ResNet18 | 400 | 79.22 | 95.06 |
| W-MSE | ResNet18 | 400 | 67.60 | 90.94 |
| Barlow Twins ++ | ResNet18 | 400 | 80.38 | 95.28 |
| NNCLR ++ | ResNet18 | 400 | 79.80 | 95.28 |
| **Musketeer** (Ours) | ResNet18 | 400 | **81.93** ← | **96.23** ← |

# Quantitative Results: Objective importance

| Method | NNCLR | Centroid | Redundancy | EMA | Acc@1 | NN Acc@1 |
|---|---|---|---|---|---|---|
| Musketeer (v0) | ✓ | ✗ | ✗ | ✗ | 69.62 | 68.8 |
| Musketeer (v1) | ✗ | ✓ | ✗ | ✗ | 67.4 | 82.8 |
| Musketeer (v2) | ✓ | ✓ | ✗ | ✗ | 71.02 | 85.28 |
| Musketeer (v3) | ✓ | ✓ | ✗ | ✓ | 71.08 ◄ | **86.16** ◄ |
| Musketeer (v4) | ✗ | ✓ | ✓ | ✓ | 71.31 | 80.6 |
| Musketeer (v5) | ✓ | ✗ | ✓ | ✓ | 71.64 | 78.8 |
| **Musketeer (v6)** | ✓ | ✓ | ✓ | ✓ | **72.17** ◄ | 82.16 |

✓ = Included  ✗ = Not included

# Qualitative Analysis: UMAP

# Qualitative Analysis: Silhouette



Overall and per class Silhouette coefficients

| | NNCLR | Musketeer |
|---|---|---|
| Overall Silhouette | 0.0006 | **0.0137** |

baby −0.0096
girl +0.0151
boy +0.0171
woman +0.0163
man +0.0069
wolf_sil +0.0206
tiger_sil +0.0237
forest −0.0054
palm +0.0009
maple +0.0008
willow −0.0046
lobster +0.0208
dolphin −0.0047
whale +0.0023
Overall +0.0131

# Qualitative Analysis: NN Retrieval

# Features of Musketeer

- Not very sensitive regarding the number of neighbours extracted.

- More expensive than single neighbour contrast.

# Speaking about Food Applications

# Food recognition



Food category and class recognition

Try it: www.logmeal.es/demo

21:22

LogMeal is a HealthApp and API in the cloud that is able to automatically recognize and analyze food from images.
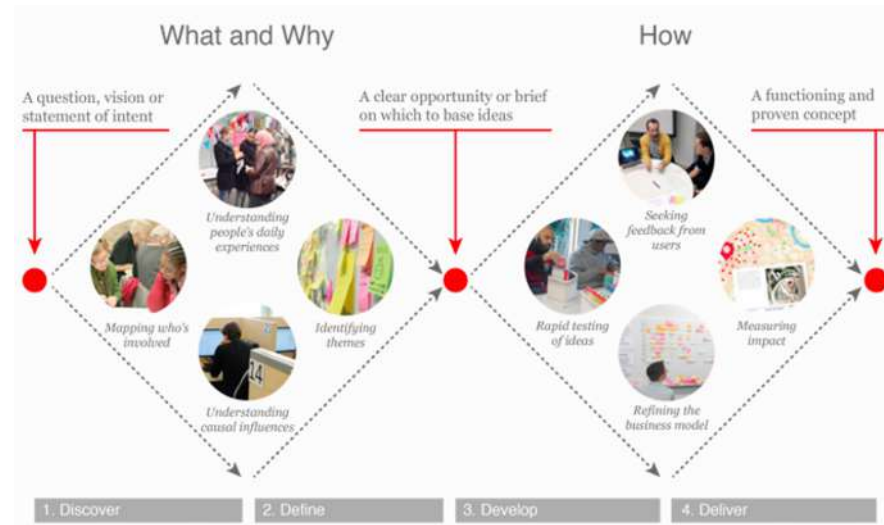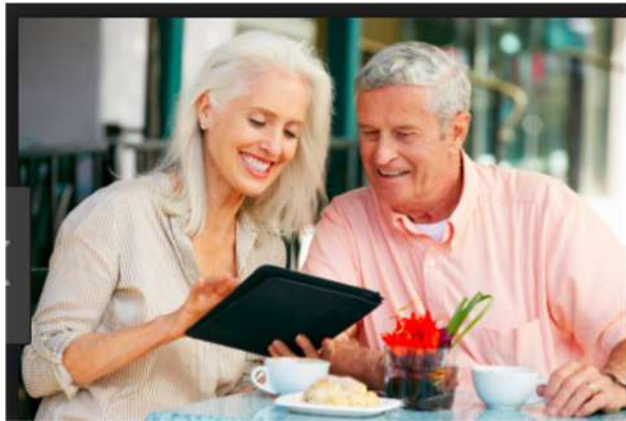
LogMeal

Validithi (EIT Health, 2019/20)

- **Automatic food diary** construction (UB).
- **Accurate, objective and continuous** food intake monitoring (UB).
- Semi-automatic **volume estimation** (Nestle).
- **Meal planner** and **health recommendations (**Nestle**).**

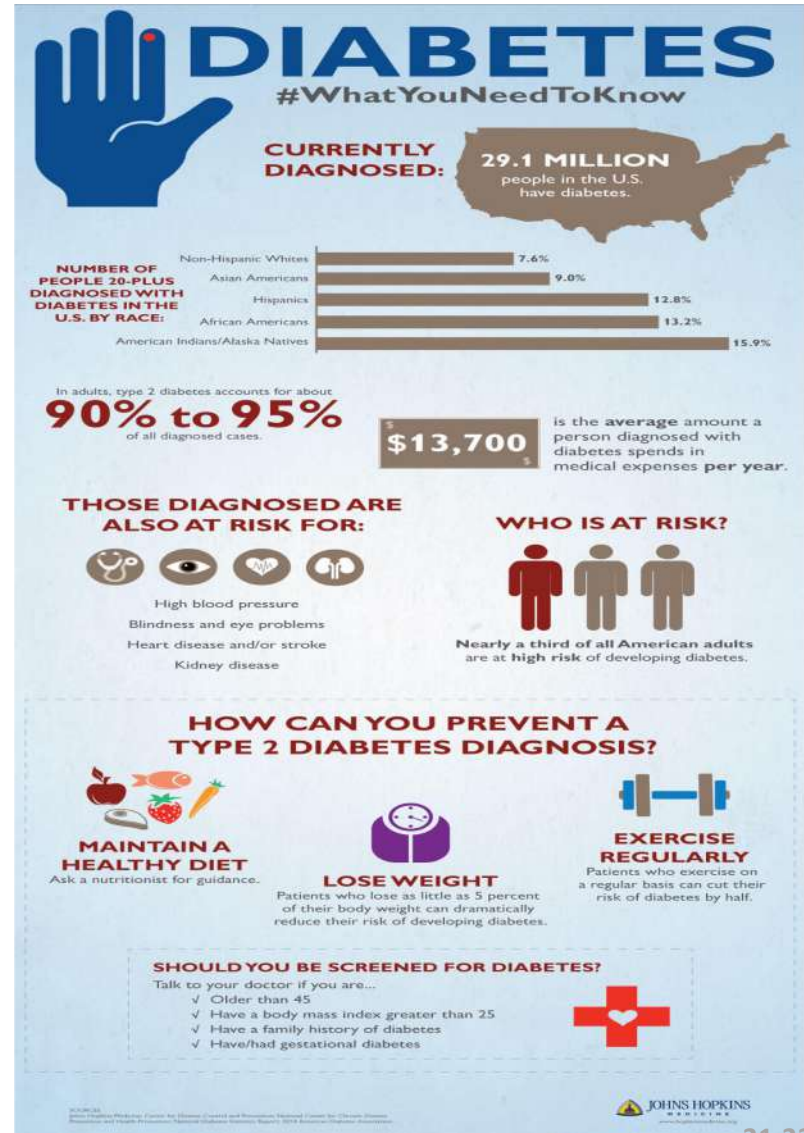| What and Why | | How | |
|---|---|---|---|
| A question, vision or statement of intent | A clear opportunity or brief on which to base ideas | | A functioning and proven concept |
| Understanding people's daily experiences | Identifying themes | Seeking feedback from users | Measuring impact |
| Mapping who's involved | | Rapid testing of ideas | |
| Understanding causal influences | | Refining the business model | |
| 1. Discover | 2. Define | 3. Develop | 4. Deliver |

NESTORE developed a multi-dimensional, personalized coaching system to support healthy ageing:

1) Generating and sustaining motivation to take care of health;

1) Suggesting healthy nutrition and personalized physical and mental coach, as well as social interaction, to prevent decline and preserve wellbeing.

The number of people with diabetes has increased from 108 million in 1980 to 422 million in 2017.

Pulso Edicions and UB are developing an app oriented to diabetic people in order to monitor their food intake and receive objective and timely feedback.
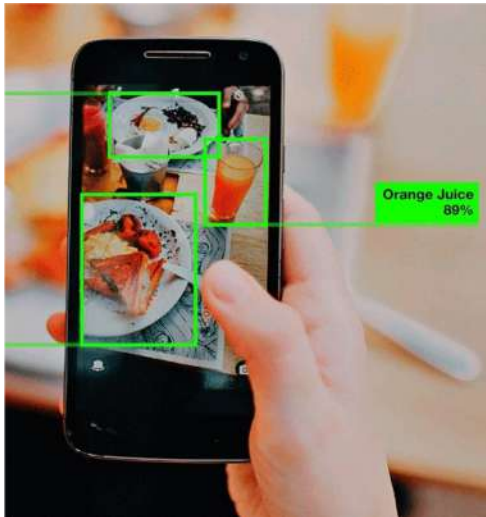




21:22

Ted Talk of Chantal Linders: "Manage the Monster in Your Head"   Greenhabit (EIT Digital, 2020/21)

21:22

# Success story: Aigecko Technologies

Touchless Checkout System: Plate and food recognition Multiple Payment System
User identification (NFC, QR, Face recognition, company card)



API that allows food recognition (ready meals and food) using Artificial Intelligence algorithms with just a photo.

21:22

# AIGecko's Food Image Analysis Applications

**FOOD TYPE DETECTION**
API to detect cooked food, prepared food, beverages, fresh vegetables and fruits, non-food products and more.

**FOOD GROUP DETECTION**
Detects the basic food groups present in food. Ideal for the generation of food records and food diaries.

**SINGLE DISH RECOGNITION**
Detects more than 880 different local and international dishes from any cuisine in the world.
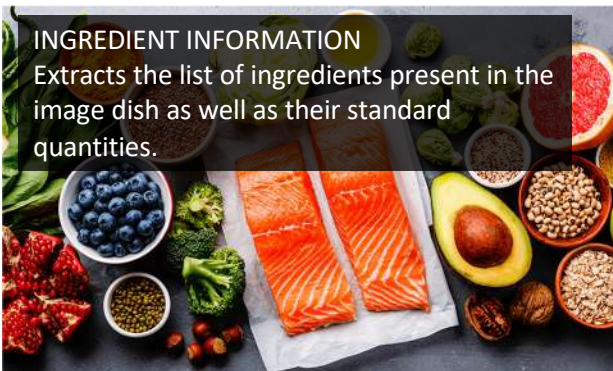
**RECOGNITION OF VARIOUS DISHES**
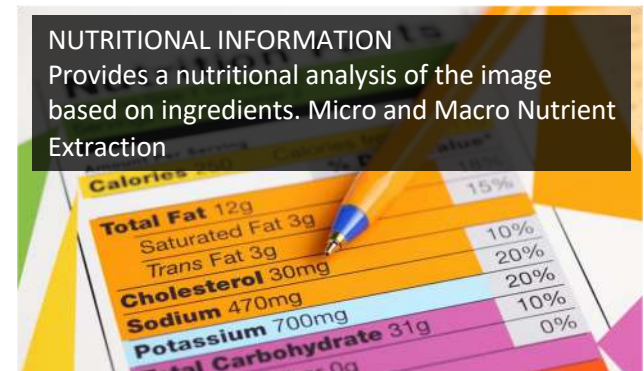Recognises and lists all the foods present in a combination dish.

**INGREDIENT INFORMATION**
Extracts the list of ingredients present in the image dish as well as their standard quantities.

**NUTRITIONAL INFORMATION**
Provides a nutritional analysis of the image based on ingredients. Micro and Macro Nutrient Extraction

21:22

# Conclusions

❖ OptSSL method **outperforms SimSiam and Barlow Twins** for the food image representation task.

   ❖ Showed the importance of **contrasting both positive and negative** samples.

❖ Musketeer introduces **Self-attention operations to create single representations, defined as centroids**, from the extracted neighbours.

   ❖ increases the neighbour retrieval accuracy while avoiding efficiency loss.

❖ Musketeer **combines its neighbour contrast objective with a feature redundancy reduction** objective, forming a symbiosis that proves to be beneficial in the overall performance of the framework.

❖ Musketeer **consistently outperforms SoTA** instance discrimination frameworks on popular image classification benchmarking datasets, namely, CIFAR-10, CIFAR-100 and ImageNet-100.

❖ **Food** Image Analysis is **highly underexplored problem** that could convert in an important **benchmark** for CV algorithms.

❖ Multiple **real applications** and **professional opportunities**

- "Pure" Reinforcement Learning (cherry)
  - The machine predicts a scalar reward given once in a while.
  - A few bits for some samples

- Supervised Learning (icing)
  - The machine predicts a category or a few numbers for each input
  - Predicting human-supplied data
  - 10→10,000 bits per sample

- Self-Supervised Learning (cake génoise)
  - The machine predicts any part of its input for any observed part.
  - Predicts future frames in videos
  - Millions of bits per sample

21:22

Thank you!