

Uczenie typu self-supervised w przetwarzaniu obrazów

ARKADIUSZ KWASIGROCH

WYDZIAŁ ELEKTROTECHNIKI I AUTOMATYKI, POLITECHNIKA GDAŃSKA

[akwasigroch.github.io](https://github.com/arkadiuszkwasiogroch)

Plan prezentacji

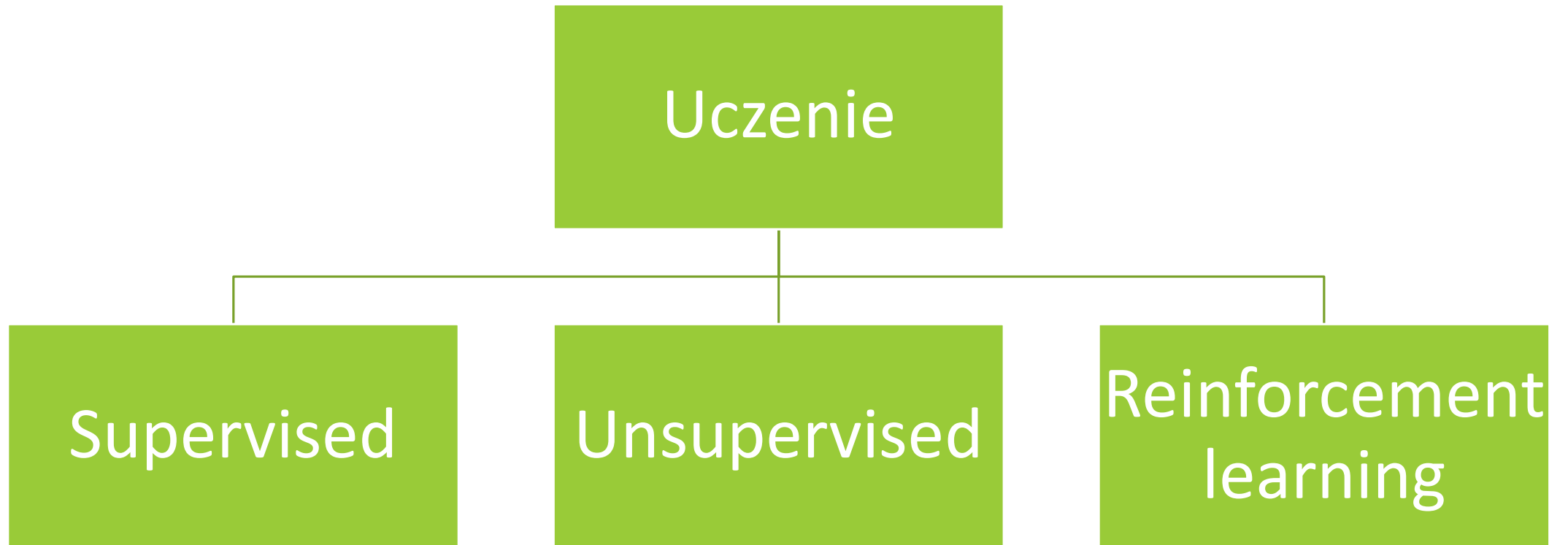
1. Typy uczenia oraz różnice pomiędzy nimi

- Supervised
- Unsupervised
- Self-supervised
- Semi-supervised
- Weakly supervised

2. Ogólny przegląd i klasyfikacja metod typu self-supervised

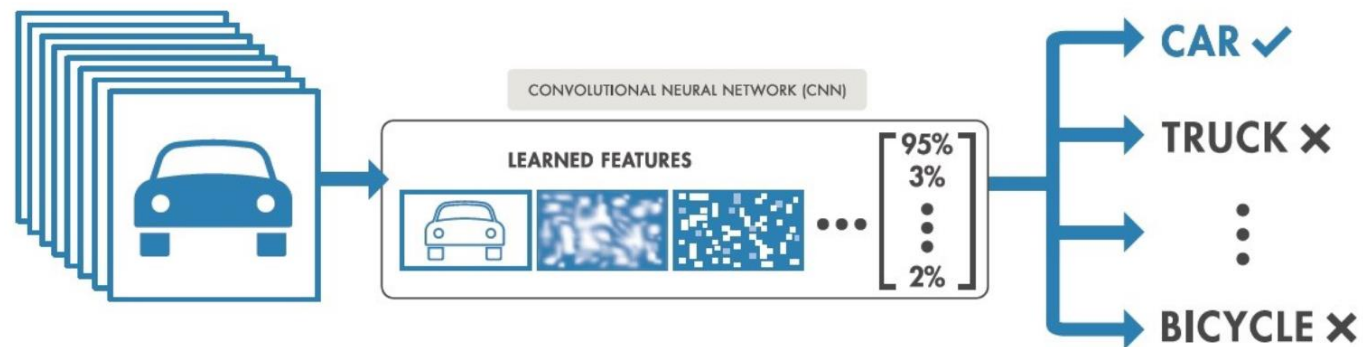
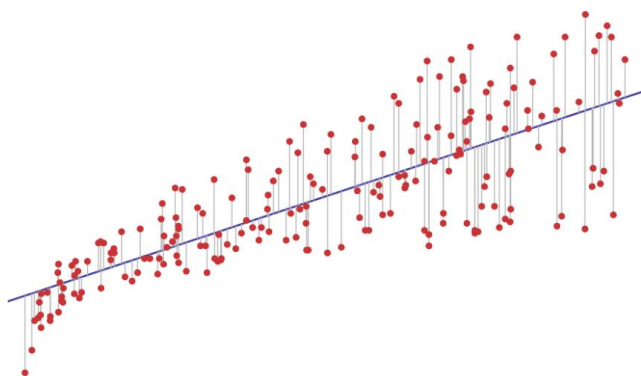
3. Szczegółowy 2 metod

Podział algorytmów



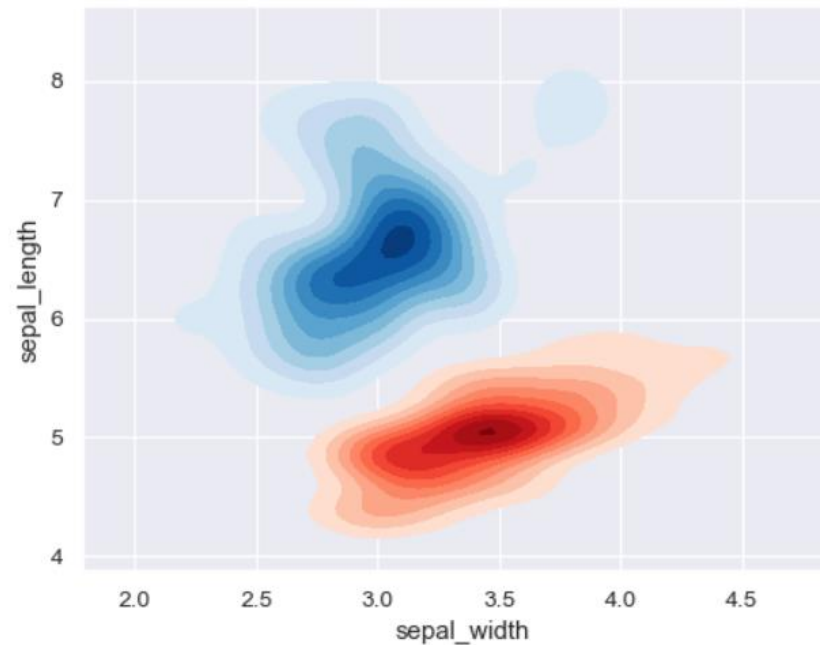
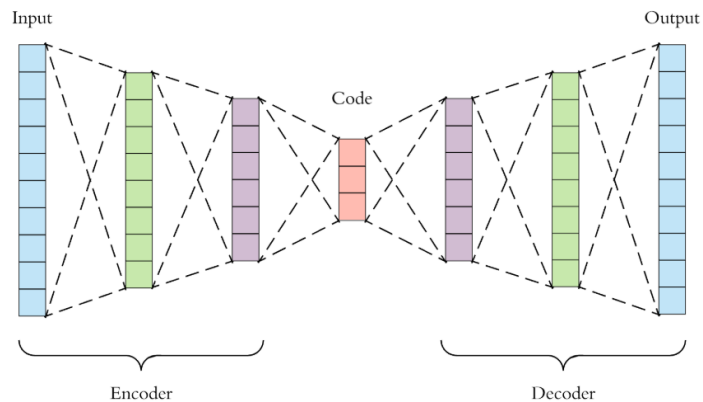
Supervised learning

- Do każdego przykładu przypisana jest etykieta
- Przykłady: problemy klasyfikacji i regresji
- Etykiety przygotowane przez człowieka (*fine-grained, human-annotated*)
- Wadą jest duży koszt przygotowania etykiet (Amazon Mechanical Turk)



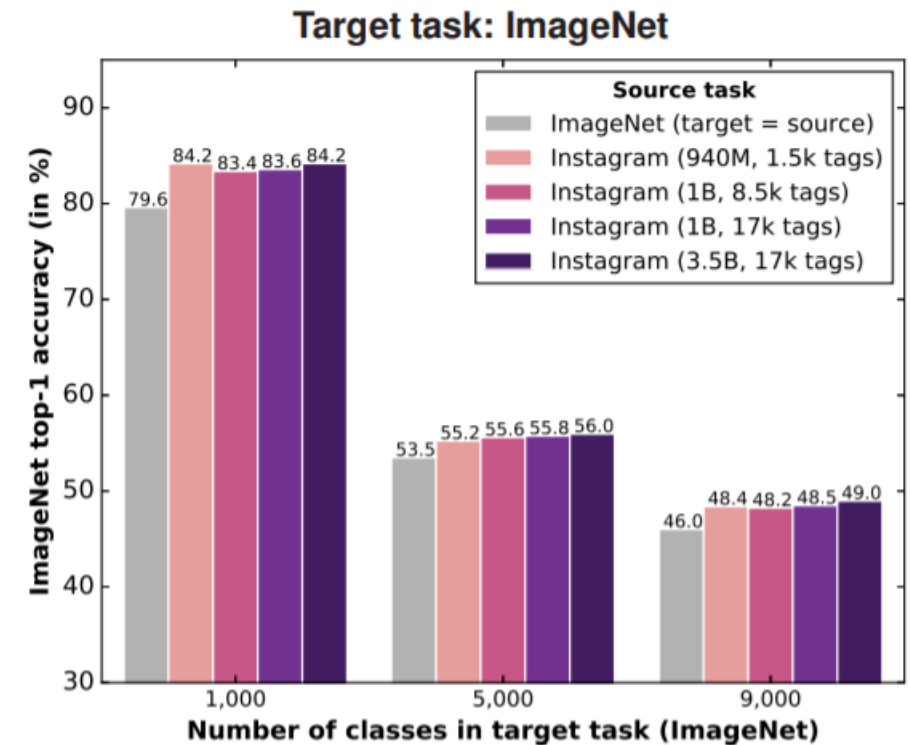
Unsupervised learning

- Algorytmy operują na danych bez etykiet
- Przykłady
 - Estymacja rozkładu
 - Klasteryzacja
 - Redukcja wymiarowości



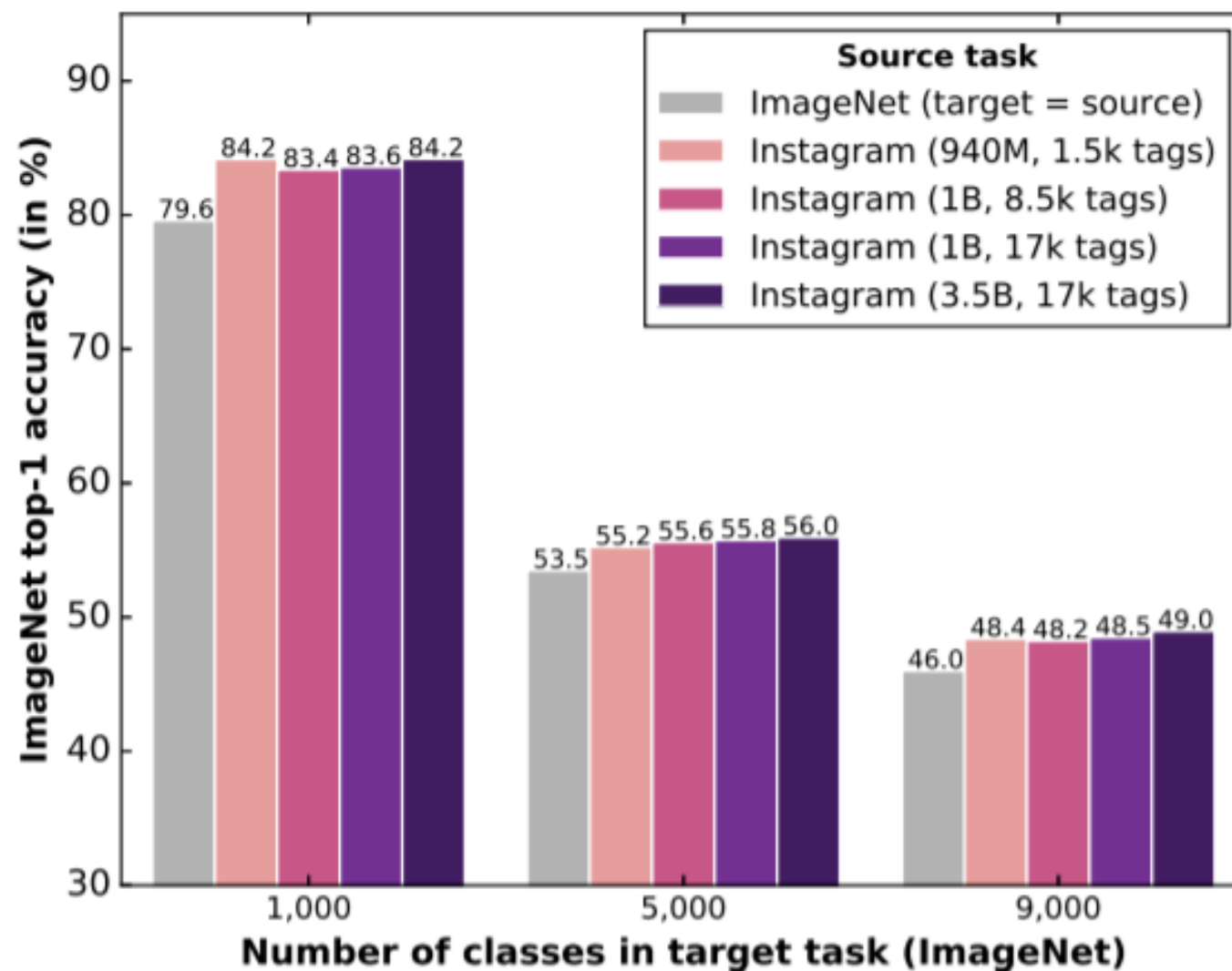
Weakly supervised

- Wykorzystanie mniej szczegółowych etykiety (coarse-grained vs fine-grained)
- Wykorzystanie etykiet niedokładnych
- Mniejszy koszt uzyskania etykiet



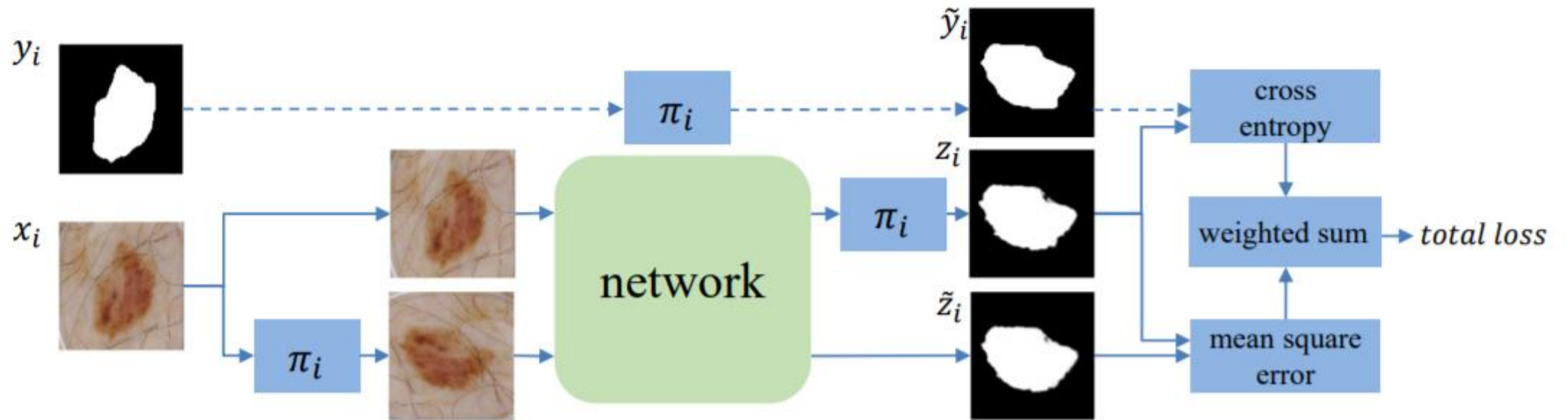
<https://arxiv.org/pdf/1805.00932.pdf>

Target task: ImageNet



Semi-supervised learning

- Włączenie do uczenia przykładów oetykietowanych i nieoetykietowanych
- Funkcja celu składa się z dwóch ważonych członów: człon supervised i człon unsupervised



Model	50 labeled, 1950 unlabeled data		
	JA	DI	SE
Supervised-only	72.85	81.15	82.77
Supervised with regularization	73.25	81.60	83.30
Our Method	75.31	83.79	86.37
Our Method-A	74.59	83.27	82.77
Our Method-B	74.21	82.68	83.15

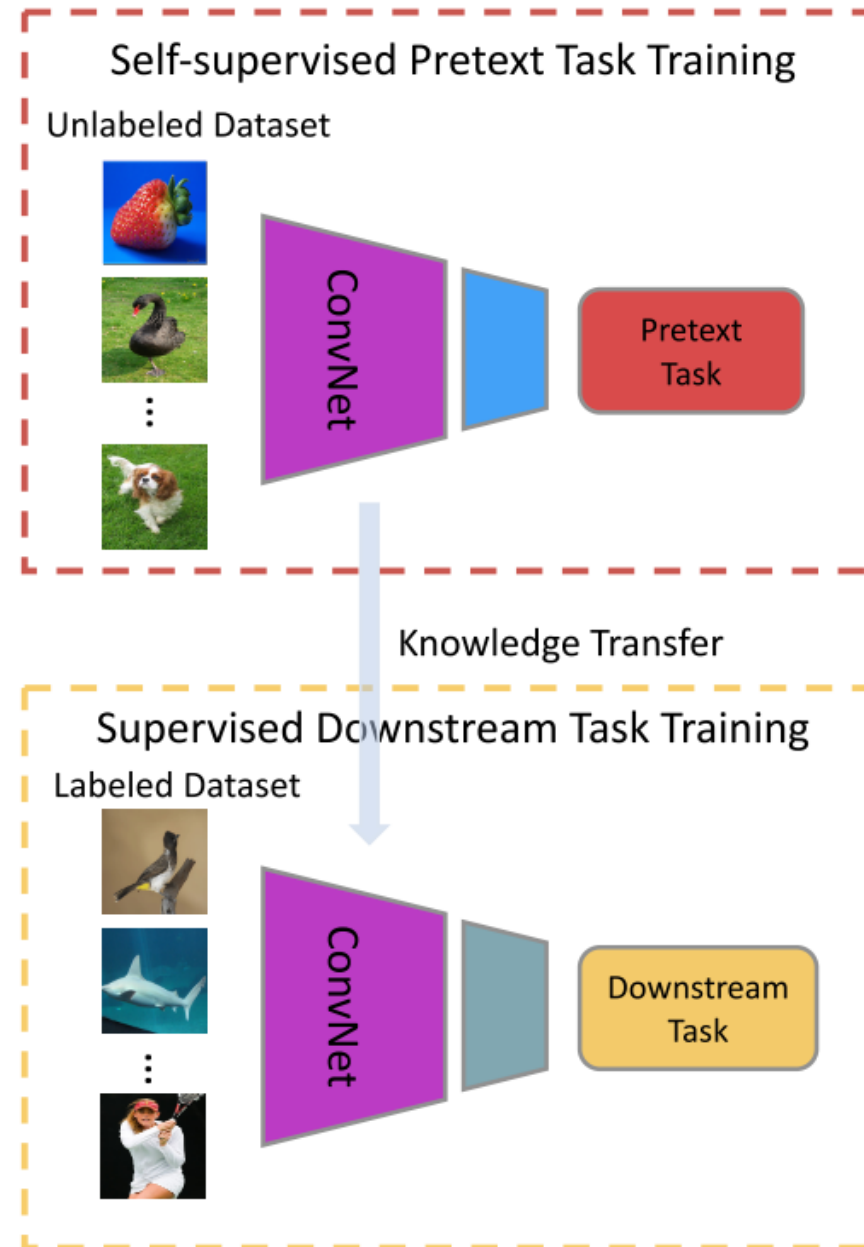
300 oetykietowanych, 1700 nieoetykietowanych
Baseline - supervised

Team	JA	DI	AC	SE	SP
Our Semi-supervised Method	0.798	0.874	0.943	0.879	0.953
Our Baseline	0.772	0.853	0.936	0.837	0.969
Yuan and Lo [28]	0.765	0.849	0.934	0.825	0.975
Berseth [8]	0.762	0.847	0.932	0.820	0.978
Bi et al. [6]	0.760	0.844	0.934	0.802	0.985
RECOD	0.754	0.839	0.931	0.817	0.970
Jer	0.752	0.837	0.930	0.813	0.976

<http://bmvc2018.org/contents/papers/0162.pdf>

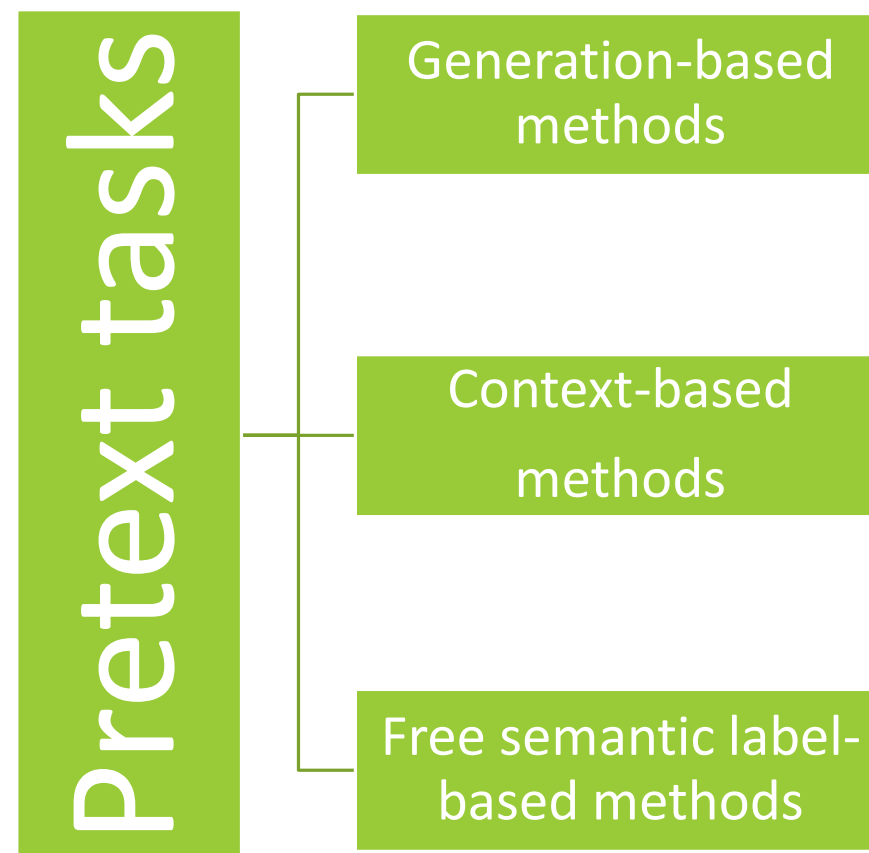
Self-supervised

- Poddziedzina unsupervised learning
- Etykiety generowane automatycznie
- Pierwszy etap – uczenie zadania pomocniczego (pretext task)
- Drugi etap – uczenie zadania docelowego (downstream task)
- Podejście podobne do transfer learning
- Zastosowanie w przetwarzaniu obrazów inspirowane rozwiązaniami NLP



Self-supervised learning – pretext tasks

- Zadanie pomocnicze (pretext task) jest zadaniem, które ma rozwiązać sieć neuronowa, w celu nauczenia odpowiedniej reprezentacji danych.
- Pseudo-etykiety generowane są automatycznie

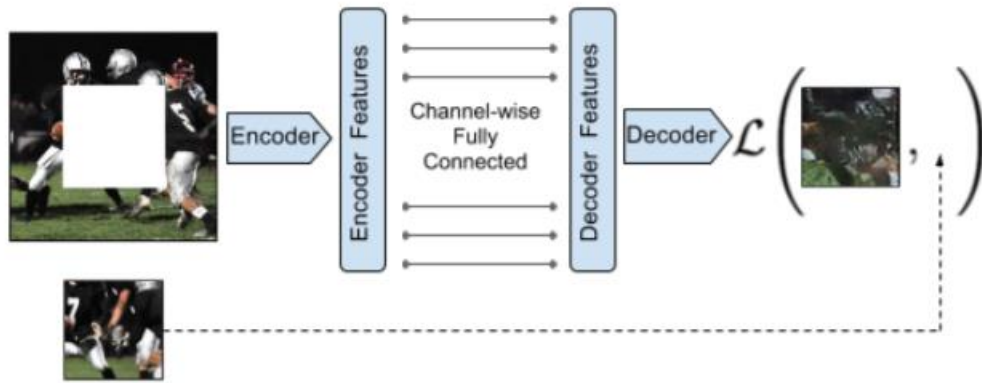


Self supervise learning – downstream task

- Zadaniem docelowym najczęściej jest klasyfikacja, detekcja lub segmentacja
- Ewaluacja przeprowadzana jest najczęściej z zamrożonymi wagami
- W przypadku klasyfikacji, wykorzystuje się klasyfikator liniowy

Generation based methods

Image inpainting



(a) Input context



(b) Human artist



(c) Context Encoder (L_2 loss)



(d) Context Encoder ($L_2 + \text{Adversarial loss}$)

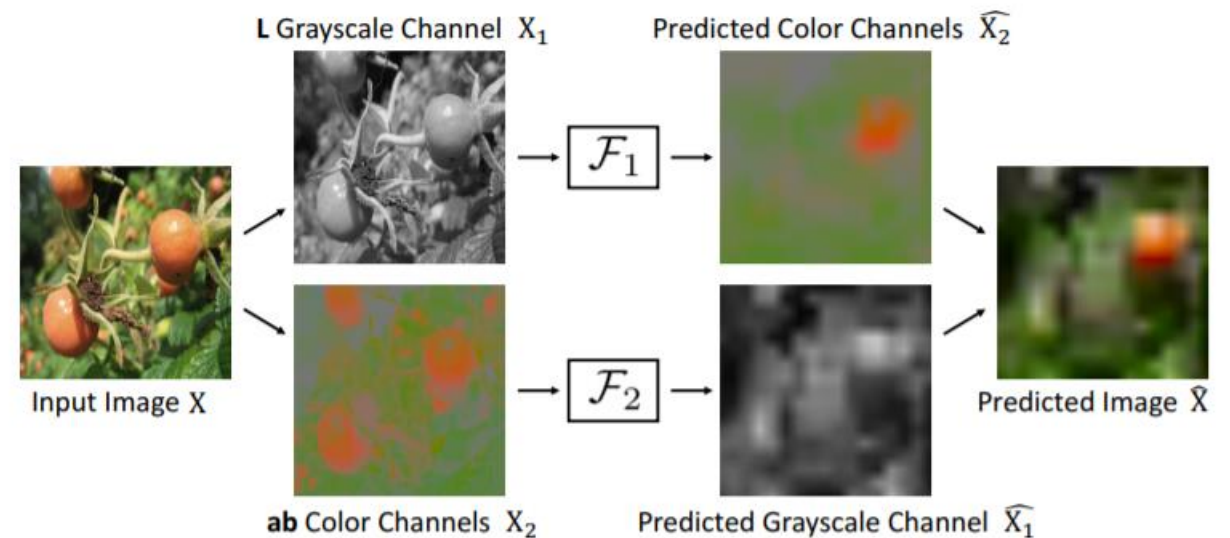
Generation based methods

Image colorization



<https://arxiv.org/pdf/1603.06668.pdf>

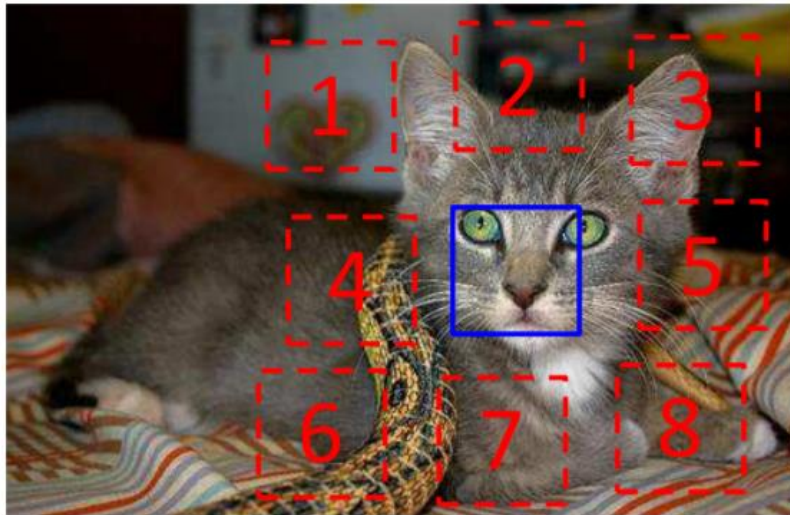
Split-Brain Autoencoders:



<https://arxiv.org/pdf/1611.09842.pdf>

Context based methods

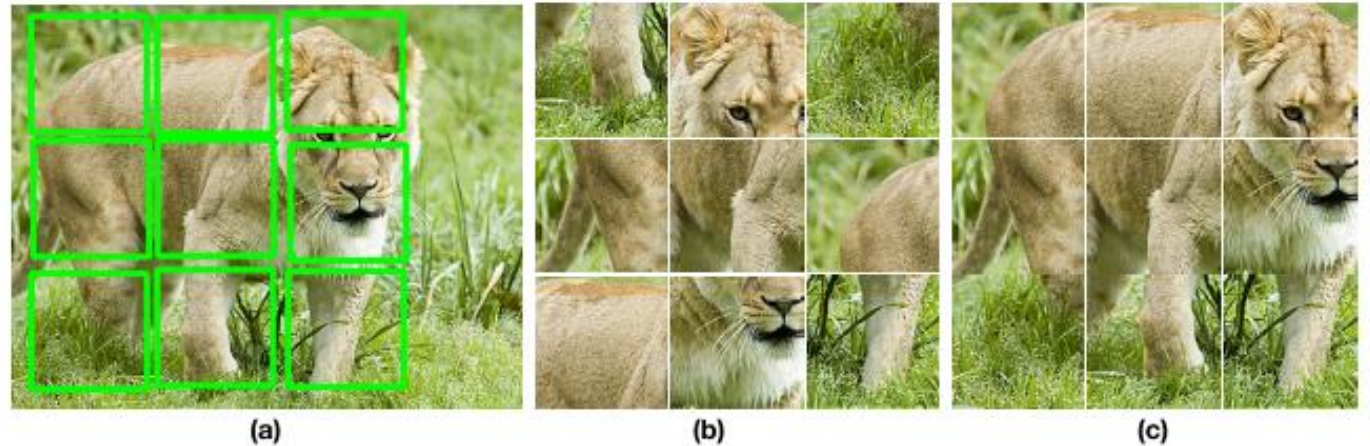
Context prediction



$$X = \left(\begin{array}{c} \text{cat face} \\ \text{cat ear} \end{array} \right); Y = 3$$

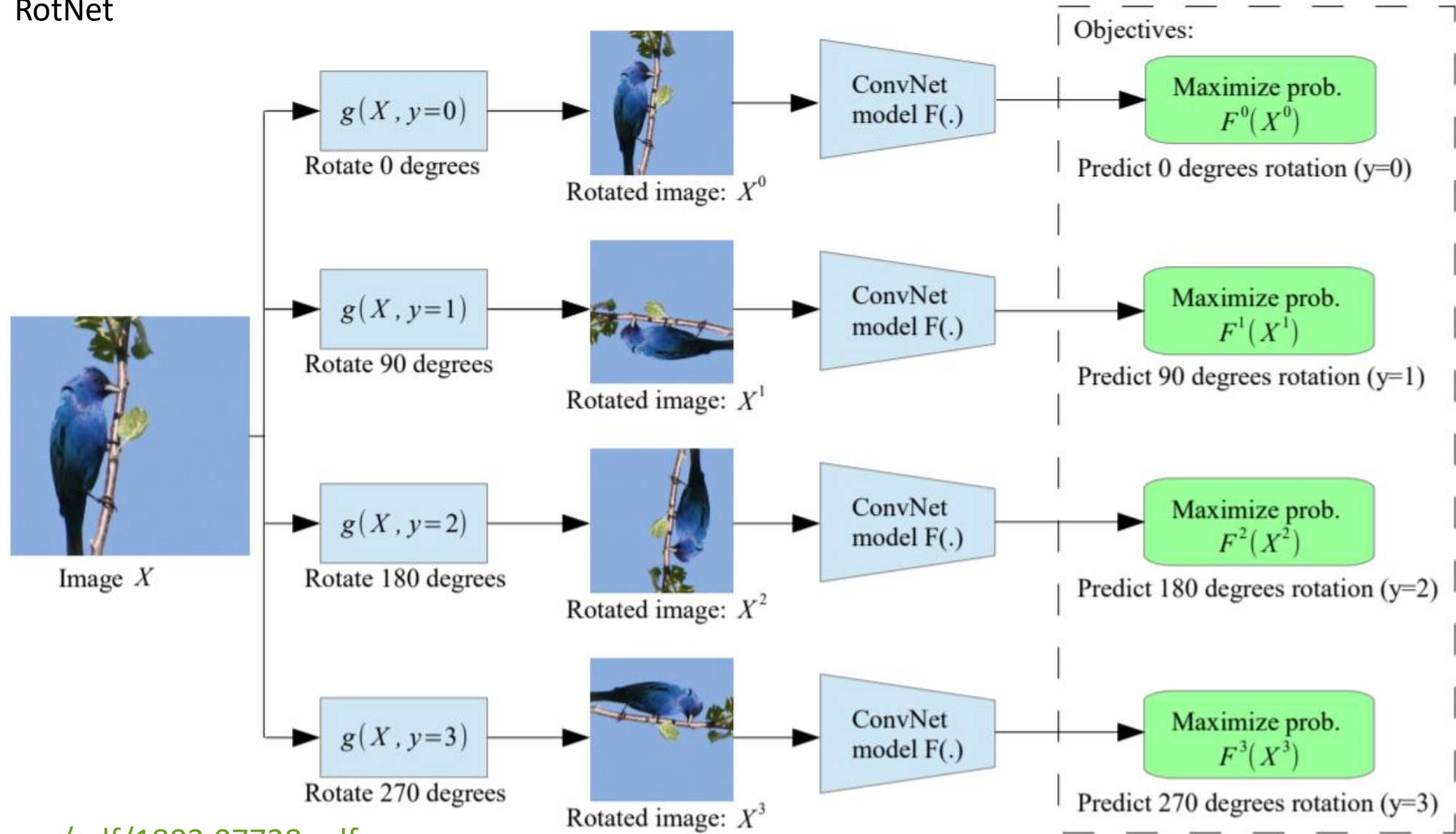
https://www.cv-foundation.org/openaccess/content_iccv_2015/papers/Doersch_Unsupervised_Visual_Representation_ICCV_2015_paper.pdf

Jigsaw puzzle solving



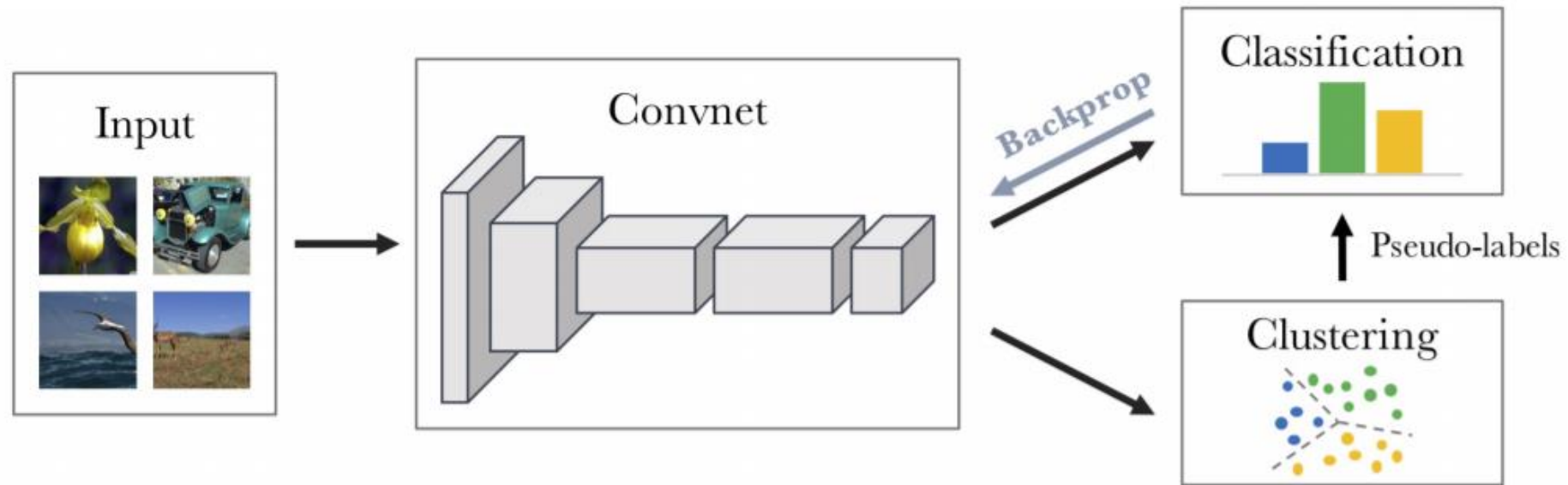
<https://arxiv.org/pdf/1603.09246.pdf>

RotNet



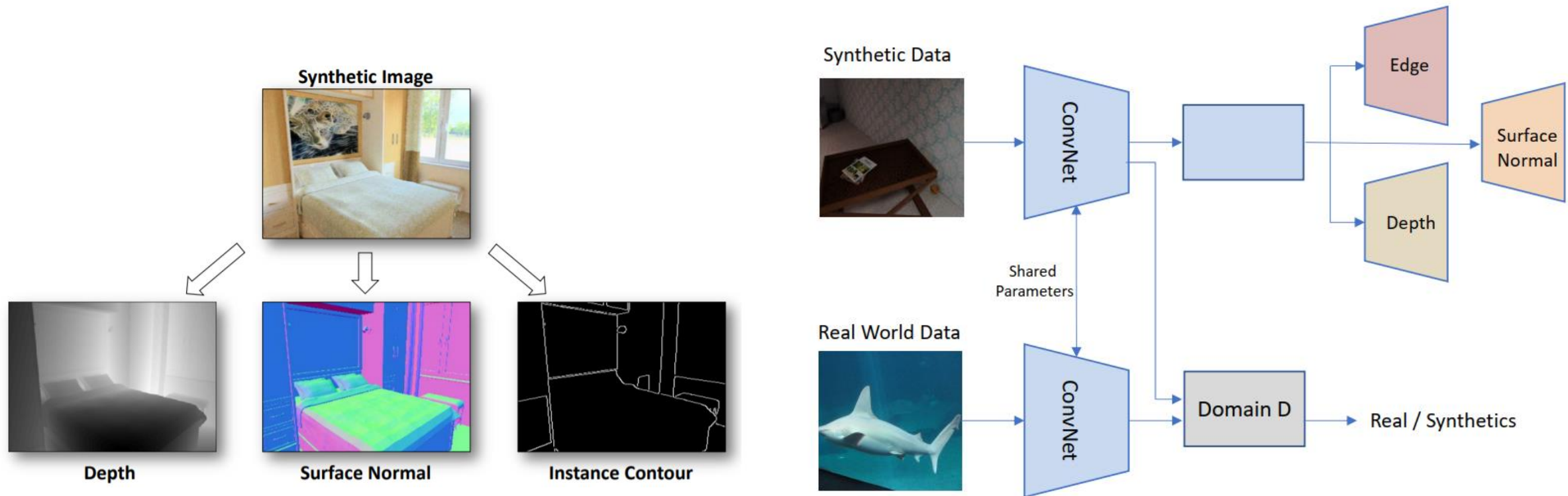
Context based methods

DeepCluster



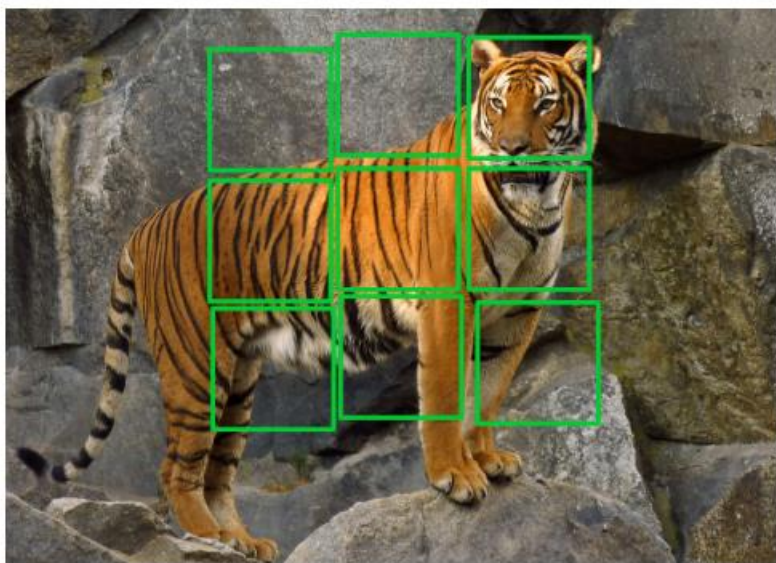
<https://arxiv.org/pdf/1807.05520.pdf>

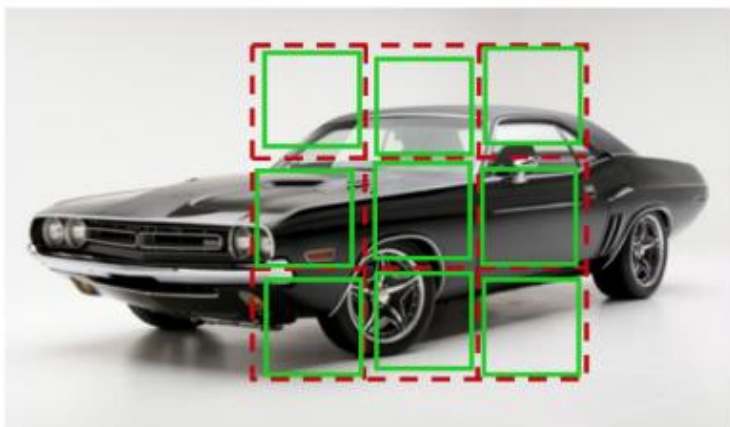
Free semantic label-based methods



<https://arxiv.org/pdf/1711.09082.pdf>

Unsupervised Learning of Visual Representations by Solving Jigsaw Puzzles

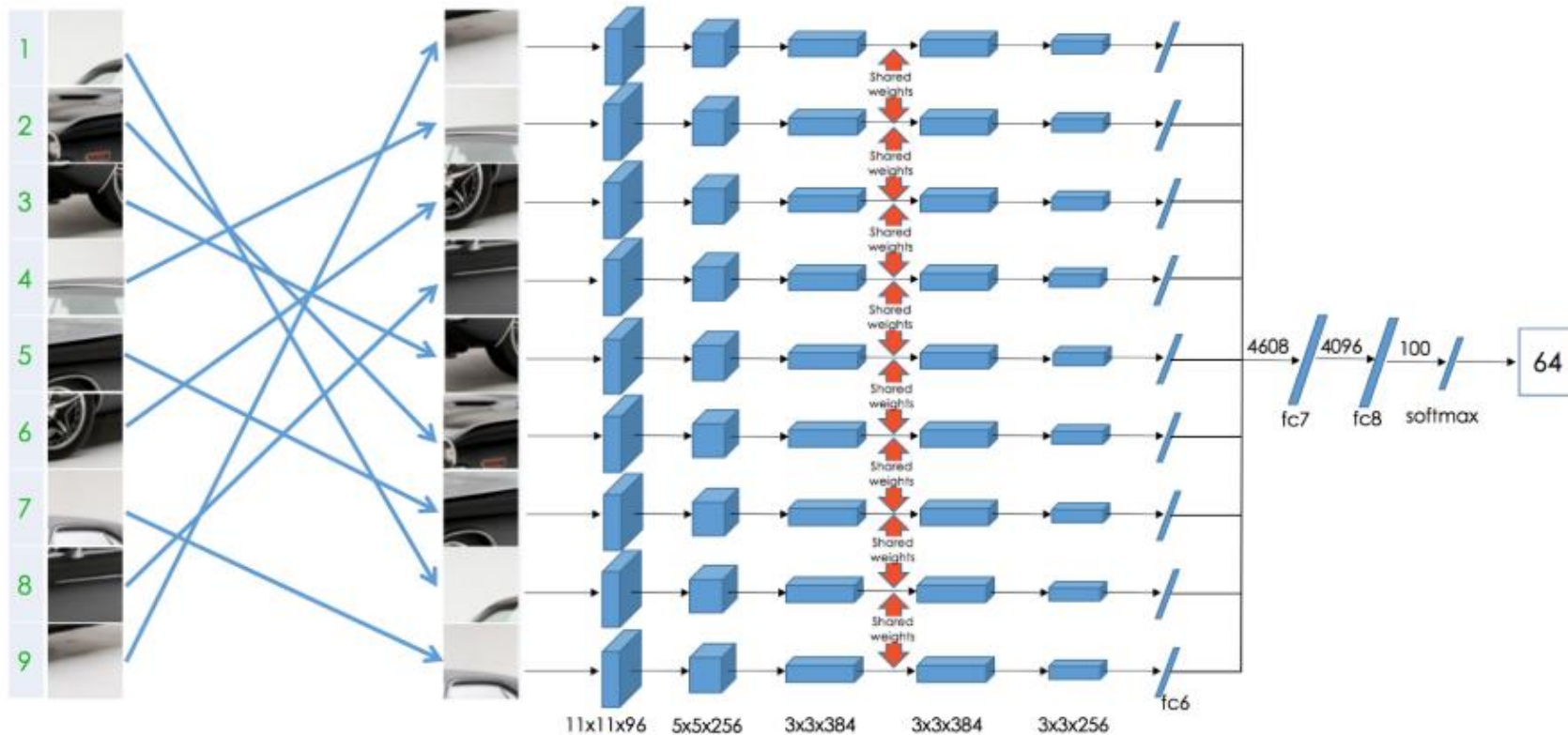




Permutation Set

index	permutation
64	9,4,6,8,3,2,5,1,7

Reorder patches according to the selected permutation



Wybór permutacji

- Liczba możliwych permutacji wynosi $9! = 362\,880$
- Należało ograniczyć liczbę permutacji i wybrać zbiór permutacji, które są do siebie „najbardziej odmienne”
- Miarą odmienności permutacji jest odległość Hamminga

Odległość Hamminga (ang. Hamming distance)– wprowadzona przez Richarda Hamminga miara odmienności dwóch ciągów o takiej samej długości, wyrażająca liczbę miejsc (pozycji), na których te dwa ciągi się różnią.

https://pl.wikipedia.org/wiki/Odleg%C5%82o%C5%9B%C4%87_Hamminga

Przykład:

4	2	3	1	8	6	7	5
4	3	2	1	8	7	6	5

Odległość Hamminga wynosi 4

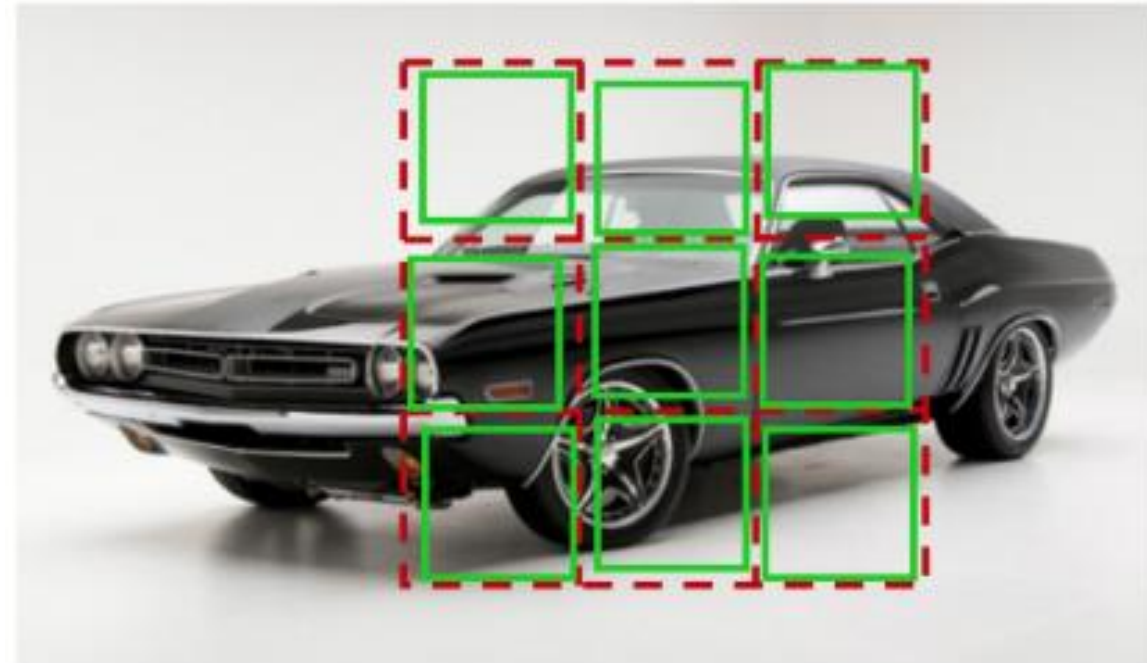
Algorithm 1. Generation of the *maximal* Hamming distance permutation set

Input: N \\ number of permutations
Output: P \\ maximal permutation set
1: $\bar{P} \leftarrow$ all permutations $[\bar{P}_1, \dots, \bar{P}_{9!}]$ \\ \bar{P} is a $9 \times 9!$ matrix
2: $P \leftarrow \emptyset$
3: $j \sim \mathcal{U}[1, 9!]$ \\ uniform sample out of $9!$ permutations
4: $i \leftarrow 1$
5: **repeat**
6: $P \leftarrow [P \ \bar{P}_j]$ \\ add permutation \bar{P}_j to P
7: $\bar{P} \leftarrow [\bar{P}_1, \dots, \bar{P}_{j-1}, \bar{P}_{j+1}, \dots]$ \\ remove \bar{P}_j from \bar{P}
8: $D \leftarrow \text{Hamming}(P, P')$ \\ D is an $i \times (9! - i)$ matrix
9: $\bar{D} \leftarrow \mathbf{1}^T D$ \\ \bar{D} is a $1 \times (9! - i)$ row vector
10: $j \leftarrow \arg \max_k \bar{D}_k$ \\ \bar{D}_k denotes the k -th entry of \bar{D}
11: $i \leftarrow i + 1$
12: **until** $i \leq N$

A good self-supervised task is neither simple nor ambiguous.

Utrudnienie zadania

- Maksymalizacja dystansu Hamminga pomiędzy permutacjami
- W celu uniknięcia uczenia na podstawie „ciągłości” pikseli elementy są od siebie oddalone (średnio o 11 pikseli)
- Każdy element normalizowany indywidualnie, wprowadzono color jittering



Wyniki

Table 1: Results on PASCAL VOC 2007 Detection and Classification. The results of the other methods are taken from Pathak *et al.* [30].

Method	Pretraining time	Supervision	Classification	Detection	Segmentation
Krizhevsky <i>et al.</i> [25]	3 days	1000 class labels	78.2%	56.8%	48.0%
Wang and Gupta [39]	1 week	motion	58.4%	44.0%	-
Doersch <i>et al.</i> [10]	4 weeks	context	55.3%	46.6%	-
Pathak <i>et al.</i> [30]	14 hours	context	56.5%	44.5%	29.7%
Ours	2.5 days	context	67.6%	53.2%	37.6%

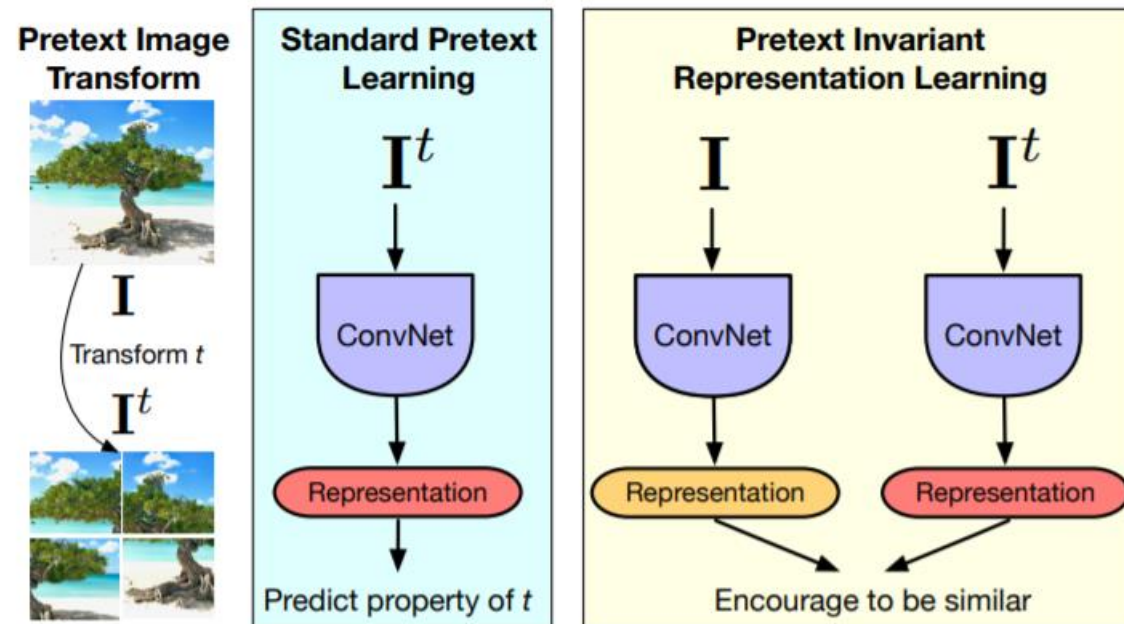
Table 4: Ablation study on the impact of the permutation set.

Number of permutations	Average hamming distance	Minimum hamming distance	Jigsaw task accuracy	Detection performance
1000	8.00	2	71	53.2
1000	6.35	2	62	51.3
1000	3.99	2	54	50.2
100	8.08	2	88	52.6
95	8.08	3	90	52.4
85	8.07	4	91	52.7
71	8.07	5	92	52.8
35	8.13	6	94	52.6
10	8.57	7	97	49.2
7	8.95	8	98	49.6
6	9	9	99	49.7

Table 5: Ablation study on the impact of the shortcuts.

Gap	Normalization	Color jittering	Jigsaw task accuracy	Detection performance
\times	\checkmark	\checkmark	98	47.7
\checkmark	\times	\checkmark	90	43.5
\checkmark	\checkmark	\times	89	51.1
\checkmark	\checkmark	\checkmark	88	52.6

Self-Supervised Learning of Pretext-Invariant Representations



<https://arxiv.org/pdf/1912.01991.pdf>

Funkcja celu

Funkcja celu zbliża do siebie cechy obrazu oryginalnego i jego modyfikacji. Oddala cechy obraz modyfikowanego oraz cechy innych obrazów ze zbioru uczącego.

$$h(\mathbf{v_I}, \mathbf{v_{I^t}}) = \frac{\exp\left(\frac{s(\mathbf{v_I}, \mathbf{v_{I^t}})}{\tau}\right)}{\exp\left(\frac{s(\mathbf{v_I}, \mathbf{v_{I^t}})}{\tau}\right) + \sum_{\mathbf{I'} \in \mathcal{D}_N} \exp\left(\frac{s(\mathbf{v_{I^t}}, \mathbf{v_{I'}})}{\tau}\right)} \quad (3)$$

$$L_{\text{NCE}}(\mathbf{I}, \mathbf{I^t}) = -\log[h(f(\mathbf{v_I}), g(\mathbf{v_{I^t}}))] \\ - \sum_{\mathbf{I'} \in \mathcal{D}_N} \log[1 - h(g(\mathbf{v_I^t}), f(\mathbf{v_{I'}}))]$$

Bank cech

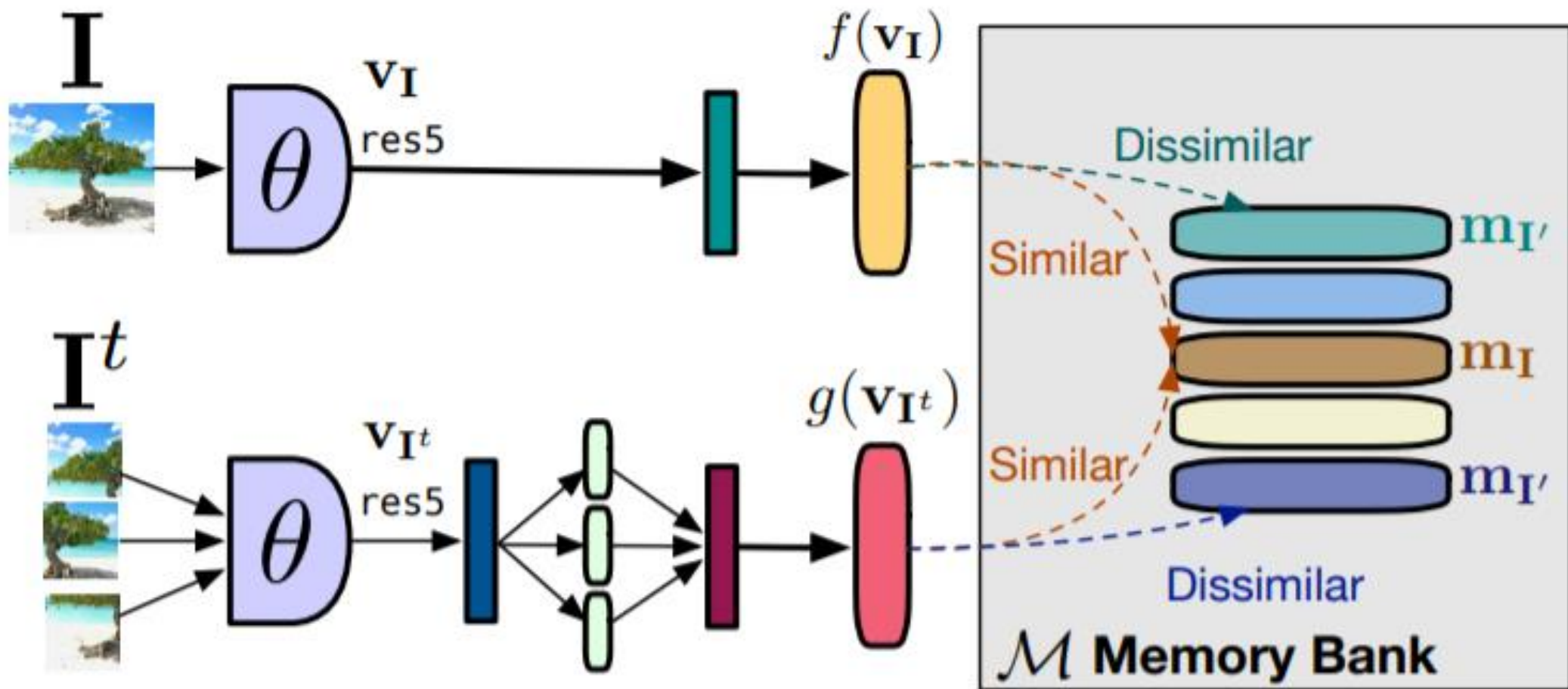
Bank cech przechowuje cechy każdego obrazu znajdującego się w zbiorze uczącym. Przechowywane reprezentacje są średnią eksplotencjalną cech oryginalnego obrazu, uzyskiwane co każdą epokę.

Nowa funkcja celu

Funkcja celu z poprzedniego slajdu nie porównywała reprezentacji obrazów oryginalnych z reprezentacją innych obrazów oryginalnych znajdujących się w bazie. W tym celu funkcję celu zmodyfikowano:

$$L(\mathbf{I}, \mathbf{I}^t) = \lambda L_{\text{NCE}}(\mathbf{m}_{\mathbf{I}}, g(\mathbf{v}_{\mathbf{I}^t})) \\ + (1 - \lambda) L_{\text{NCE}}(\mathbf{m}_{\mathbf{I}}, f(\mathbf{v}_{\mathbf{I}})).$$

- Pierwszy człon funkcją celu z poprzedniego slajdu. Zamiast aktualnie obliczanej reprezentacji, wykorzystuje się reprezentację przechowywaną w pamięci
- Drugi człon posiada dwie role:
 - porównywanie reprezentacji w pamięci z oryginalną reprezentacją zapobiegają gwałtownym zmianom parametrów
 - Oddala reprezentację oryginalnego obrazu od reprezentacji innych oryginalnych obrazów znajdujących się w bazie



Szczegóły techniczne

- Sieć ResNet50 jako ekstraktor cech
- Sieć ResNet50 zwraca wektor 2048 elementowy, stosowana jest dodatkowa warstwa w pełni połączona, w celu zredukowania tego wektora do rozmiaru 128 elementów
- Zadaniem pomocniczym (pretext task) jest „układanie puzzli”, każdy element generuje wektor 128 elementowy, następnie wektory są łączone ze sobą, tworząc wektor $9 * 128$ elementów. Taki wektor podawany jest na warstwę w pełni połączoną w celu uzyskania wektora 128 elementowego.
- Współczynnik funkcji celu 0.5

Wyniki - detekcja

Detekcja obrazu na bazie VOC07+12, pretraining na bazie Imagenet

Method	Network	AP ^{all}	AP ⁵⁰	AP ⁷⁵	Δ AP ⁷⁵
Supervised	R-50	52.6	81.1	57.4	=0.0
Jigsaw [19]	R-50	48.9	75.1	52.9	-4.5
Rotation [19]	R-50	46.3	72.5	49.3	-8.1
NPID++ [72]	R-50	52.3	79.1	56.9	-0.5
PIRL (ours)	R-50	54.0	<u>80.7</u>	59.7	+2.3

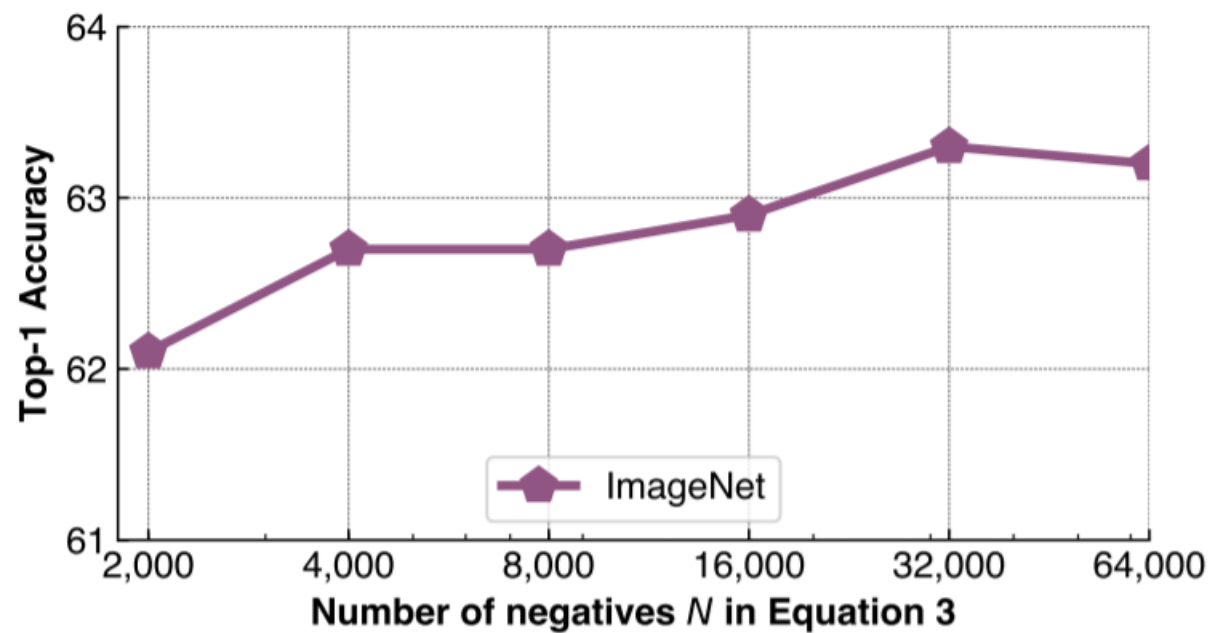
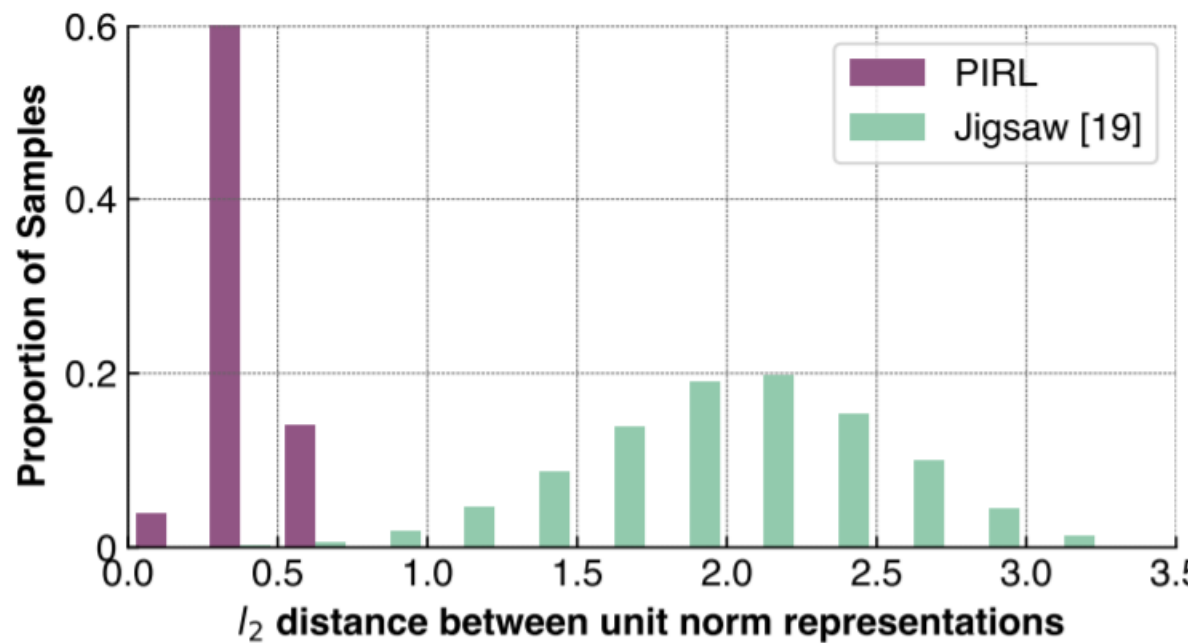
Wyniki – klasyfikacja liniowa

Method	Parameters	Transfer Dataset			
		ImageNet	VOC07	Places205	iNat.
ResNet-50 using evaluation setup of [19]					
Supervised	25.6M	75.9	87.5	51.5	45.4
Colorization [19]	25.6M	39.6	55.6	37.5	–
Rotation [18]	25.6M	48.9	63.9	41.4	23.0
NPID++ [72]	25.6M	59.0	76.6	46.4	32.4
MoCo [24]	25.6M	60.6	–	–	–
Jigsaw [19]	25.6M	45.7	64.5	41.2	21.3
PIRL (ours)	25.6M	63.6	81.1	49.8	34.1

Wyniki – klasyfikacja przy małych zbiorach uczących

Method	Data fraction →	1%	10%
	Backbone	Top-5 Accuracy	
Random initialization [72]	R-50	22.0	59.0
NPID [72]	R-50	39.2	77.4
Jigsaw [19]	R-50	45.3	79.3
NPID++ [72]	R-50	52.6	81.5
VAT + Ent Min. [20, 45]	R-50v2	47.0	83.4
S ⁴ L Exemplar [75]	R-50v2	47.0	83.7
S ⁴ L Rotation [75]	R-50v2	53.4	83.8
PIRL (ours)	R-50	57.2	83.8

Wyniki



Wyniki – inne zadania pomocnicze

Method	Params	Transfer Dataset			
		ImageNet	VOC07	Places205	iNat.
Rotation [18]	25.6M	48.9	63.9	41.4	23.0
PIRL (Rotation; ours)	25.6M	60.2	77.1	47.6	31.2
Δ of PIRL	-	+11.3	+13.2	+6.2	+8.2
Combining pretext tasks using PIRL					
PIRL (Jigsaw; ours)	25.6M	62.2	79.8	48.5	31.2
PIRL (Rotation + Jigsaw; ours)	25.6M	63.1	80.3	49.7	33.6

Prezentacja przygotowana na podstawie przeglądu:

<https://arxiv.org/pdf/1902.06162.pdf>