

# Cdiscount's Image Classification

MinWoo Byeon



# About me

- MinWoo Byeon ([mwbyun@hotmail.com](mailto:mwbyun@hotmail.com))
- Software Engineer at [Kakao Corp](#)
- Lives in [Jeju Island](#), South Korea
- Bachelor's degree, Computer Engineering
- Interested in Machine Learning, Algorithms and Competitions
- Kaggle: <https://www.kaggle.com/mwbyeon>
- LinkedIn: <https://www.linkedin.com/in/minwoo-byeon/>

*(I'm not fluent in English. Please let me know if I need to fix it.)*

# Kaggle

## Welcome to Kaggle Competitions

Challenge yourself with real-world machine learning problems

---



### New to Data Science?

Get started with a tutorial on our most popular competition for beginners, [Titanic: Machine Learning from Disaster](#).



### Build a Model

Get the data & use whatever tools or methods you prefer to make predictions.

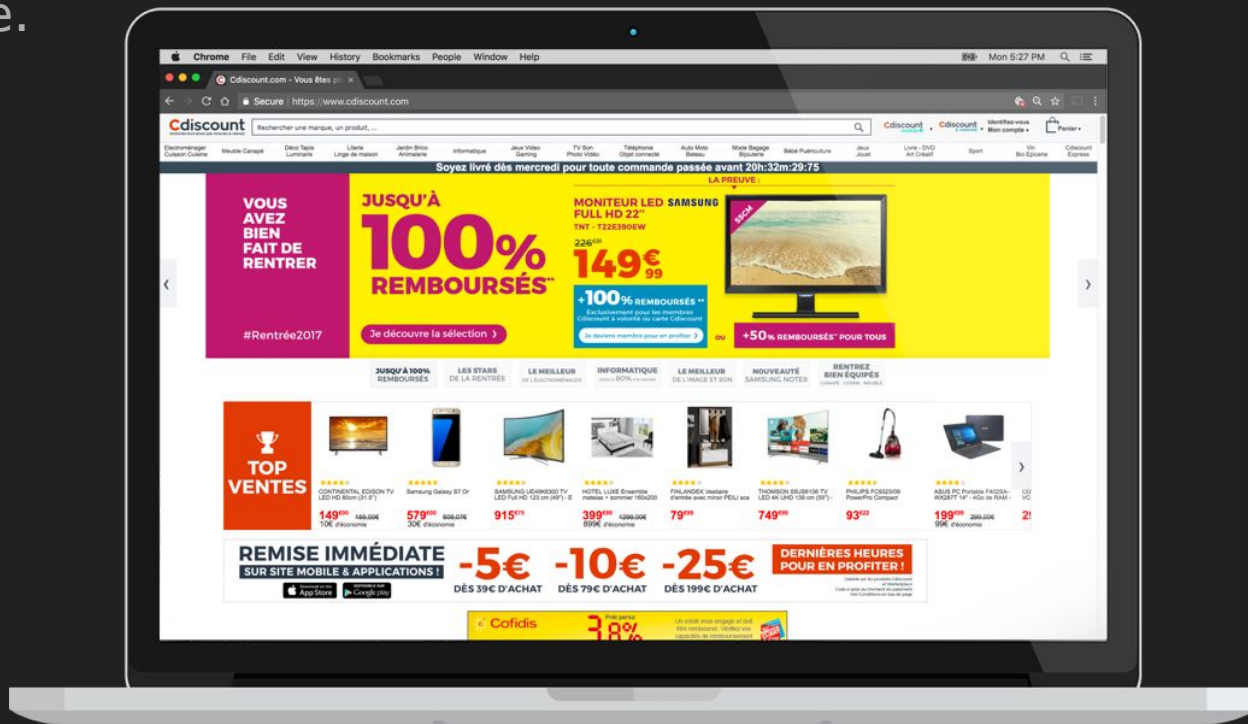


### Make a Submission

Upload your prediction file for real-time scoring & a spot on the leaderboard.

# Cdiscount's Image Classification Challenge

- Cdiscount.com is the largest non-food e-commerce company in France.



# Cdiscount's Image Classification Challenge

- <https://www.kaggle.com/c/cdiscount-image-classification-challenge>
  - In this challenge you will be building a model that **automatically classifies the products based on their images**. As a quick tour of Cdiscount.com's website can confirm, one product can have one or several images. The data set Cdiscount.com is making available is unique and characterized by superlative numbers in several ways:
    - Almost **9 million products**: half of the current catalogue
    - Each product contains between **1-4 images**.
    - More than **15 million images** at **180x180** resolution
    - More than **5200 categories**: yes this is quite an extreme multi-class classification!

# Competition Environments

- Python3 (with PyCharm)
- MXNet 0.12.0 (GPU enabled)
- Pascal based 8-GPU Machine

Dataset

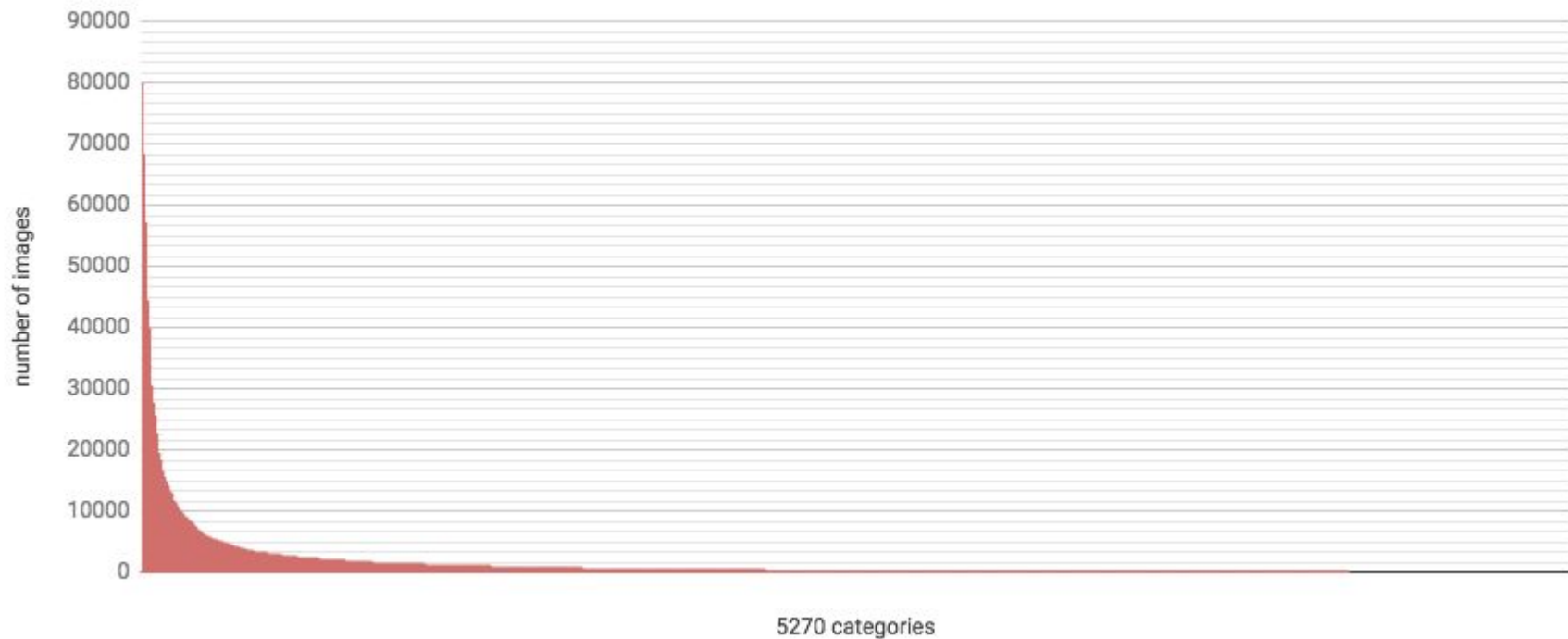
# Dataset

- 5270 Categories
- Training Set
  - train.bson, 58.2 GB
  - 7,069,896 products
  - 12,371,293 images
    - each product contains between 1-4 images
- Test Set
  - test.bson, 14.5GB
  - 1,768,182 products



# Imbalanced Dataset

Category Distribution



# Training and Validation

- Split the BSON file into Training and Validation
- Split products in the `train.bson`  
into `train_train.bson` and `train_valid.bson`
  - randomly selected with seed
  - Training(0.95) : Validation(0.05)

# Create different datasets

- DATASET A
  - split products into 0.95(training) : 0.05(validation)
  - random seed: 12648430 (0xC0FFEE)
- DATASET B
  - split products into 0.95(training) : 0.05(validation)
  - random seed: 1
- DATASET C
  - split products into 0.95(training) : 0.05(validation)
  - random seed: 12648430 (0xC0FFEE)
  - remove duplicated images from DATASET A (it can reduce training time)

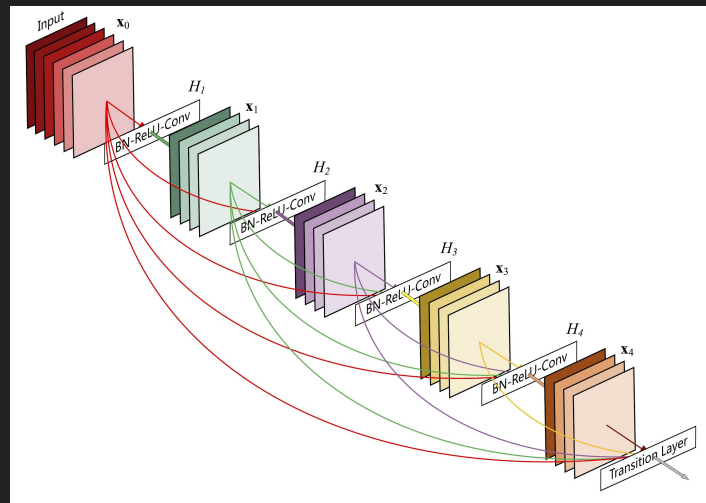
# Convolutional Neural Networks

# CNNs

- DenseNets
  - <https://arxiv.org/abs/1608.06993>
- ResNext
  - <https://arxiv.org/abs/1611.05431>
- DPNs(Dual Path Networks)
  - <https://arxiv.org/abs/1707.01629>
- SE-ResNext
  - <https://arxiv.org/abs/1709.01507>

# DenseNets

- Paper
  - Densely Connected Convolutional Networks
  - <https://arxiv.org/abs/1608.06993>
- Pre-trained Models using MXNet
  - <https://github.com/liuzhuang13/DenseNet>
  - DenseNet-121
  - DenseNet-169
  - DenseNet-201
  - DenseNet-161



# ResNext

- Paper
  - Aggregated Residual Transformations for Deep Neural Networks
  - <https://arxiv.org/abs/1611.05431>
- Pre-trained Models using MXNet
  - <https://github.com/apache/incubator-mxnet/tree/master/example>
  - <http://data.mxnet.io/models/imagenet/resnext/>
  - ResNext-50
  - ResNext-101
  - ResNext-101-64x4d

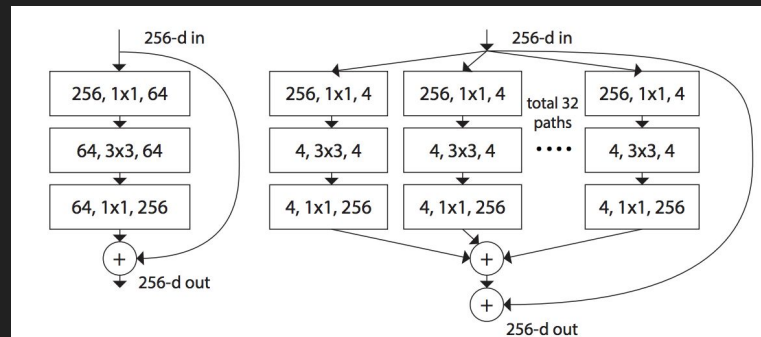
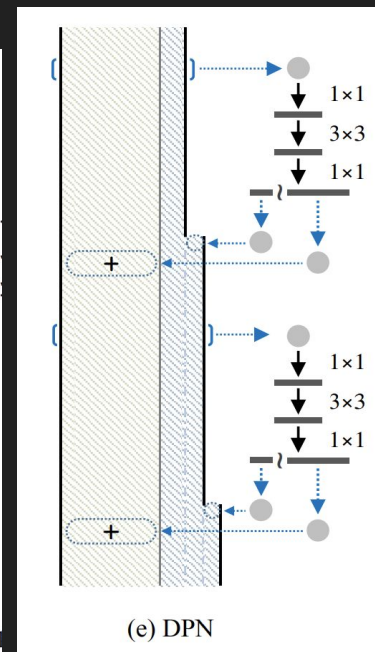
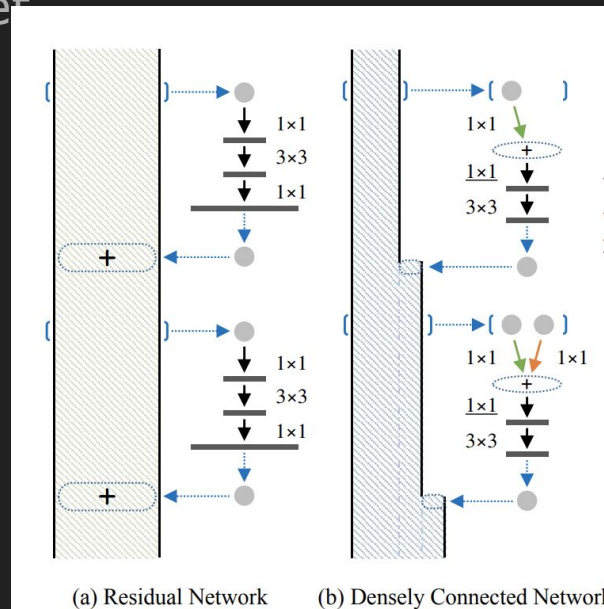


Figure 1. **Left:** A block of ResNet [14]. **Right:** A block of ResNeXt with cardinality = 32, with roughly the same complexity. A layer is shown as (# in channels, filter size, # out channels).

# DPNs(Dual Path Networks)

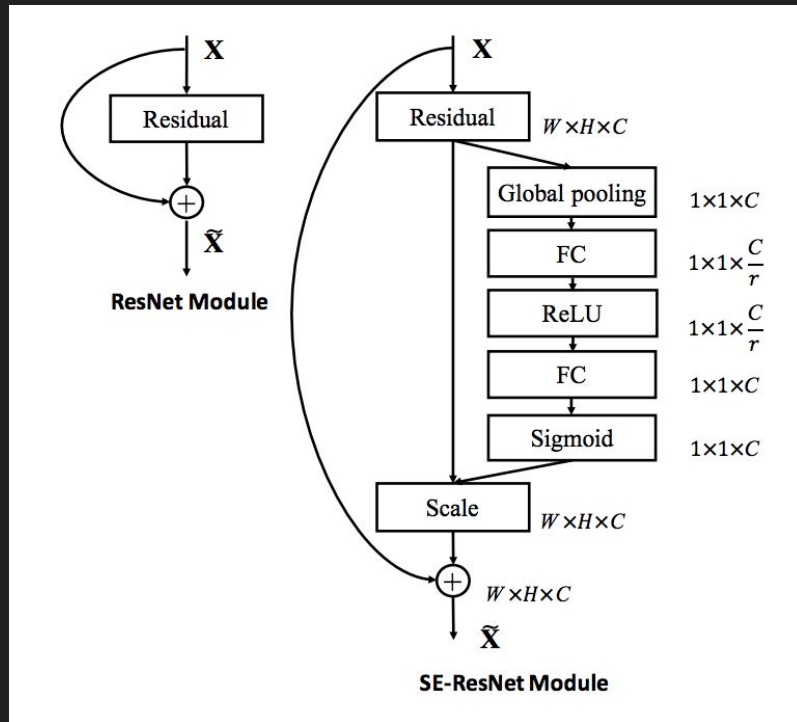
- Paper
  - Dual Path Networks
  - <https://arxiv.org/abs/1707.01629>
- Pre-trained Models using MXNet
  - <https://github.com/cypw/DPNs>
  - DPN-68
  - DPN-92
  - DPN-98
  - DPN-131
  - DPN-107





# SE-ResNext

- Paper
  - Squeeze-and-Excitation Networks
  - <https://arxiv.org/abs/1709.01507>
- Pre-trained Models using MXNet
  - Not Exists
  - Transfer some parameters from ResNext
- Implementation
  - I modified the ResNext and implemented it myself.  
(see `train/symbol/resnext.py`)



# Training

# Parameter Initialization

- Transfer Learning from ImageNet Models
  - used pre-trained models as much as possible.
- MSRA PReLU Initializer
  - This initializer is proposed for initialization related to ReLu activation, it makes some changes on top of **Xavier method**.
  - it used last layer for classification
  - <https://mxnet.incubator.apache.org/api/python/optimization/optimization.html#mxnet.initializer.MSRAPrelu>
  - Introduced in <https://arxiv.org/abs/1502.01852>

# Optimizer

- SGD
- NAG
- Adam
- **NAdam** → It's the best in this competition.
  - Adam with Nestorov momentum

# Augmentation

- input: 180x180x3
- did not use random crop
- use only random flip (it's enough for training on 15~20 epochs)
- experiment
  - ResNet-34, Batch=512, Opt=Adam, Base-LR=0.0005, Epoch=15

Input Size	Option	Local Train (Top-1)	Local Validation (Top-1)
180 x 180	No Aug.	0.862802	0.688090
<b>180 x 180</b>	<b>Flip</b>	<b>0.787349</b>	<b>0.692491</b>
180 x 180	Flip + HSL	0.774842	0.690422
224 x 224 (NN)	Flip	<b>0.789405</b>	<b>0.695116</b>

# Image Size

- The larger the image size, the higher the accuracy.
  - used 180x180 because the larger the image size, the slower the learning.
- Experiment
  - ResNet-34, Batch=512, Opt=Adam, Base-LR=0.0005, Epoch=15

Input Size	conv0	stage1	stage2	stage3	stage4	Train-Acc	Val-Acc	Training Speed (img/sec)
160 x 160	80 x 80	40 x 40	20 x 20	10 x 10	5 x 5	0.753975	0.678730	1130/sec
180 x 180	90 x 90	45 x 45	23 x 23	12 x 12	6 x 6	0.750308	0.680361	872/sec
192 x 192	96 x 96	48 x 48	24 x 24	12 x 12	6 x 6	<b>0.751713</b>	0.680706	884/sec
224 x 224	112 x 112	56 x 56	28 x 28	14 x 14	7 x 7	0.748252	<b>0.681723</b>	647/sec

# Label Smoothing

- <https://arxiv.org/abs/1512.00567> (see Section. 7)
- Label smoothing(0.1) increases accuracy.
- Experiment
  - ResNet-34, Batch=512, Opt=Adam, Base-LR=0.0005, Epoch=15

Smooth alpha	Local Train (Top-1)	Local Validation (Top-1)
0.0	<b>0.785019</b>	0.691917
<b>0.1</b>	<b>0.779468</b>	<b>0.694035</b>
0.2	0.770593	0.692766

# GAP + Dropout

- Did not use dropout after GAP(Global Average Pooling) layer in such as ResNext, SE-ResNext
- Experiment
  - ResNet-34, Batch=512, Opt=NAdam, Base-LR=0.0005, Epoch=15

Dropout Ratio	Local Train (Top-1)	Local Validation (Top-1)
<b>0.0</b>	<b>0.788686</b>	<b>0.693604</b>
0.2	0.743764	0.689954



# Trained 14 Models

	Network	Transfer from	Dataset	Epochs	Single Val-acc
M01	ResNext-101	ImageNet-1K	C	23	0.660329
M02	DPNs-92	ImageNet-1K	C	13	0.662091
M03	DPNs-92	ImageNet-1K	C	13	0.663739
M04	DenseNet-161	ImageNet-1K	A	20	0.726930
M05	DPNs-98	ImageNet-1K	A	20	0.732641
M06	DPNs-92	ImageNet-1K	A	11	0.734462
M07	ResNext-101-64x4d	ImageNet-1K	A	15	0.735808

# Trained 14 Models

	Network	Transfer from	Dataset	Epochs	Single Val-acc
M08	DPNs-131	ImageNet-1K	A	18	0.736697
M09	ResNext-101-64x4d	ImageNet-1K	B	15	0.737427
M10	ResNext-101	ImageNet-1K	A	19	0.738542
M11	SE-ResNext-101-64x4d	M10	A	15	0.739272
M12	DPNs-131	ImageNet-1K	B	17	0.742388
M13	SE-ResNext-101-64x4d	M09	B	13	0.743221
M14	DPNs-107	ImageNet-1K	A	20	0.743781

# Training Speeds

	Batch Size	GPUs	Memory Usage (per GPU)	Speed (8-GPU, Sync)	Epoch Time
<b>ResNext-101-32x4d</b>	512	8	12,461MiB	579 images / sec	5.9 Hours
<b>ResNext-101-64x4d</b>	512	8	18,185 MiB	338 images / sec	10.2 Hours
<b>SE-ResNext-101-32x4d</b>	512	8	14,057 MiB	411 images / sec	8.3 Hours
<b>SE-ResNext-101-64x4d</b>	512	8	19,805 MiB	274 images / sec	12.5 Hours
<b>DPNs-92</b>	512	8	-	530 images /sec	6.5 Hours
<b>DPNs-98</b>	512	8	16,367 MiB	380 images / sec	9 Hours
<b>DPNs-107</b>	<b>256</b>	<b>8</b>	<b>17,707 MiB</b>	<b>246 images / sec</b>	<b>14 Hours</b>
DenseNet-161	512	8	-	570 images /sec	6 Hours

Prediction

# Testing Time Augmentation

- use original image and flipped image  
more augmentation(crop, flop, transpose) reduces accuracy.
- compute the arithmetic mean for the probability of a image

# Ensembles

- Ensemble of images in a product
  - Compute the arithmetic mean (sum of probabilities of each image)
- Ensemble of models
  - Compute the arithmetic mean (sum of probabilities of each model)
  - The geometric mean outperform the arithmetic mean in less than 6 model ensembles.
- number of inferences for single image
  - M: number of models for ensemble
  - K: TTA(Testing Time Augmentation)
  - It needs  $(M * K)$  inferences to compute the probability of single image
- number of inferences for Test Data
  - $14 \text{ models} * 2 \text{ augmentation} * 2 \text{ images(average)} * 1,768,182 = 100\text{M} : ($

# MD5 Hash

- Create a dictionary of MD5 Hash
  - compute the MD5 hash of each image in the training set.
  - compute the MD5 hash of each product in the training set.  
the hash of the product is a set of md5 hash of each image in the same product.
  - create a **dictionary of (hash, category)** pairs
  - same image/product may have multiple category. (confused...)

# MD5 Hash

- Predict a probability of images using MD5 dictionary
  - during inference time,  
if the hash of a image exists in the dictionary,  
the probability of a image is set to one-hot vector.
  - after using this method,
    - $0.77100 \rightarrow 0.77429$  (+0.00329, Private LB)
    - $0.75890 \rightarrow 0.76259$  (+0.00369, Private LB)
    - $0.75604 \rightarrow 0.76083$  (+0.00479, Private LB)



# MD5 Hash

- if same image/product have multiple category → confused...
  - See `predict/predict.py#L133-L161` for details
- Experiment
  - ResNext-101-64x4d (epoch 15, original+flip)

Mode	val-acc
None	0.758480
Unique	0.761724
Majority	0.757597
L1	0.761747
L2	0.761574
Softmax	0.753390

# MD5 Hash

- Post-processing to improve accuracy using MD5 hash of products
  - after inference time,  
if the hash of a product exists in the dictionary,  
category of this product is overwritten by the category in the dictionary.
  - after using this method,
    - 0.78998 -> 0.79046 (+0.00048, Private LB)  
(it's BEST SCORE for me)




















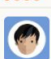
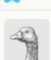
# Accelerate Pre-processing

- needs multi-processing for fast pre-processing  
but, `multiprocessing.Queue` module is very slow.
- I use [ZeroMQ](#) instead of `multiprocessing.Queue` for process communication.

# Leaderboard with Ensembles

Ensembles	Models	Private LB	Public LB
1	M14	0.76713	0.76544
2	M5, M7	0.77429	0.77282
4	M11~M14	0.78693	0.78490
8	M06~M14	0.78775	0.78577
11	M04~M14	0.78879	0.78673
14	M01~M14	0.78998	0.78813
14	M01~M14 (post-processing)	0.79046	0.78861
22	M01~M14 + extras	0.78739	0.78556

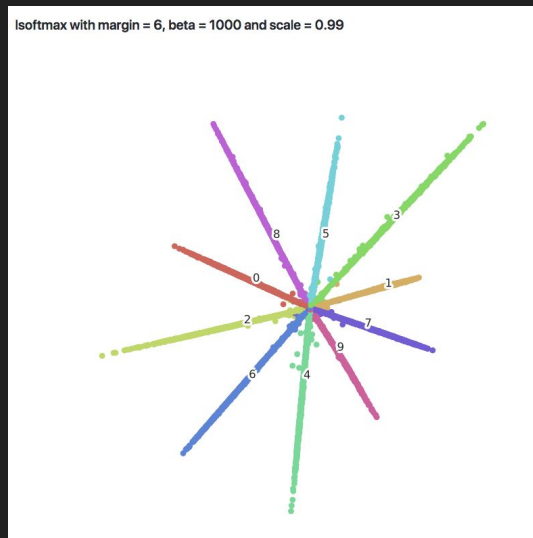
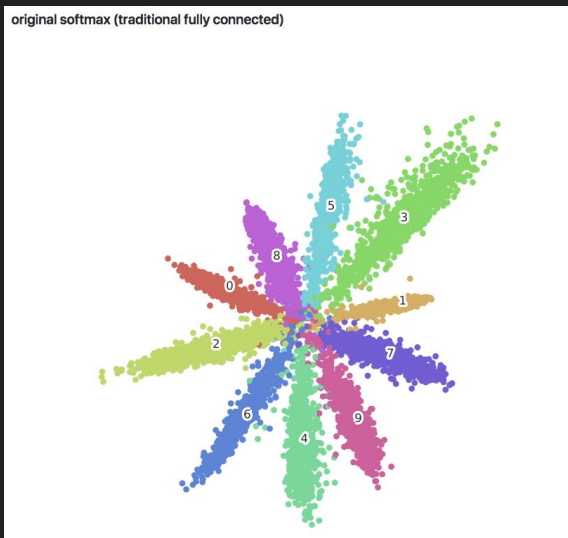
# 3rd place

#	△pub	Team Name	Kernel	Team Members	Score ?	Entries
1	—	<b>bestfitting</b>			0.79567	63
2	—	<b>Convolutud predictions</b>		  	0.79352	155
3	—	<b>Dylan @ Kakao MMRDT</b>			0.79046	90
4	—	<b>Dimensionality Diabolists</b>			0.78868	95
5	—	<b>10011000</b>		    	0.78693	169
6	—	<b>smlyaka</b>		 	0.78315	96
7	—	<b>DeepTortoise</b>		   	0.78178	77
8	—	<b>Azat Davletshin</b>			0.77883	12
9	—	<b>n01z3</b>			0.77735	51
10	▲ 1	<b>owruby</b>			0.77611	61
11	▼ 1	<b>Guanshuo Xu</b>			0.77604	6

# Other Methods

# Large Margin Softmax

- <https://arxiv.org/abs/1612.02295>
- MXNet implementation: <https://github.com/luoyetx/mx-lsoftmax>
- training failed... :(



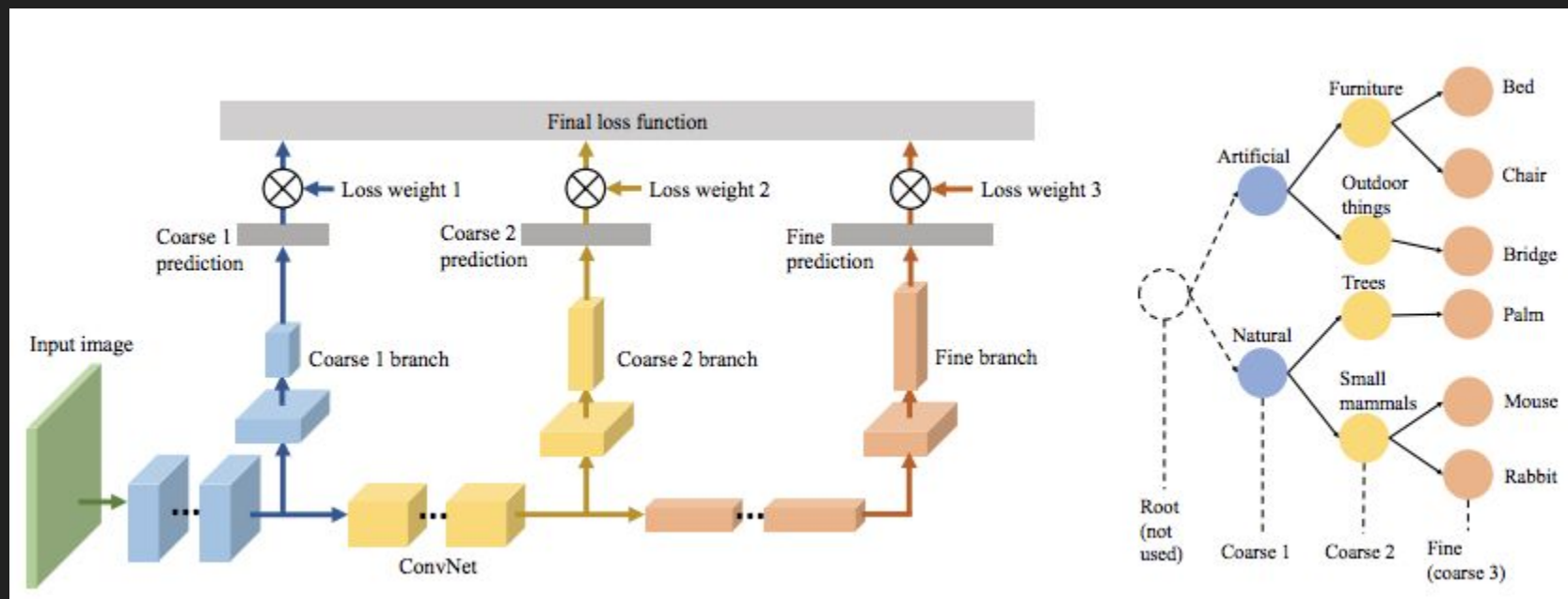
# OCR

- CD, Book
- Tried Tesseract OCR to improve accuracy for Book and Album
  - <https://github.com/tesseract-ocr/tesseract>
  - I tested some sample images
  - but, The OCR results were not accurate



# B-CNN

- <https://arxiv.org/abs/1709.09890>
  - "a good initialization may downgrade the benefit of B-CNN"



# 12-Channel CNN

- Combine four images into a 12-channel image.
- I did not have enough time for training :(

# ILSVRC 2017 Classification Competition

- A. The number of first  $1 \times 1$  convolutional channels for each bottleneck building block was halved to reduce the computation cost of the network with a minimal decrease in performance
- B. The first  $7 \times 7$  convolutional layer was replaced with three consecutive  $3 \times 3$  convolutional layers.
- C. The down-sampling projection  $1 \times 1$  with stride-2 convolution was replaced with a  $3 \times 3$  stride-2 convolution to preserve information

# ILSVRC 2017 Classification Competition

- D. A dropout layer (with a drop ratio of 0.2) was inserted before the classifier layer to prevent overfitting
- E. Label-smoothing regularisation was used during training.
- F. The parameters of all BN layers were frozen for the last few training epochs to ensure consistency between training and testing

# Source Code

<https://github.com/mwbyeon/kaggle-cdiscount>

# References

- ILSVRC 2017
  - [http://image-net.org/challenges/beyond\\_ilsvrc](http://image-net.org/challenges/beyond_ilsvrc)
  - SE-Nets: [http://image-net.org/challenges/talks\\_2017/SENet.pdf](http://image-net.org/challenges/talks_2017/SENet.pdf)
  - DPNs: [http://image-net.org/challenges/talks\\_2017/ilsvrc2017\\_DPNs.pdf](http://image-net.org/challenges/talks_2017/ilsvrc2017_DPNs.pdf)
  - Short: [http://image-net.org/challenges/talks\\_2017/ilsvrc2017\\_short\(poster\).pdf](http://image-net.org/challenges/talks_2017/ilsvrc2017_short(poster).pdf)
- ILSVRC 2016
  - Hikvision ImageNet 2016:  
[http://image-net.org/challenges/talks/2016/Hikvision\\_at\\_ImageNet\\_2016.pdf](http://image-net.org/challenges/talks/2016/Hikvision_at_ImageNet_2016.pdf)

# Appendix: 1st Place Solution

# Team: bestfitting

- Fine-tuning pre-trained models
  - 0.759/0.757, inception-resnet-v2
  - 0.757/0.756, ResNet-50
- Make full use of multi-images of a product
  - Split the train-set into 4 parts, and concatenate the K images into a image.
  - 0.772/0.771 Inception-Resnet-v2
  - 0.769/0.768+ ResNet-50 models
- Use OCR to add semantics to the models
  - use CTPN to extract the boxes containing text,the result was quite good
  - extracted the features of boxes from CRNN and feed them into my multi-input CNN
- Ensemble
  - DenseNet-161, DenseNet-169, DPN-92, ResNet-50, Inception-ResNet-v2



# Appendix: 2nd Place Solution

# Team: Convoluted predictions

- Ensemble

