



# Know Your Bands: A Content-Based Image Retrieval Mobile App

Final Presentation

14252392

SUN Jingxuan

Supervisor: Prof. Yuen Pongchi

Observer: Prof. Xu Jianliang

2018.4.21

# Outline

- Introduction
  - Motivation
  - Functionality
  - Pipeline
- Sem 1: Building pipeline
  - Data Collection
  - Object Recognition
  - Feature Extraction
  - Instance Matching
  - Information Retrieval
  - Initial test and problems
- Sem 2: Optimization
  - Network optimization
  - Feature fusion experiments
- Mobile development
- Q&A

# Introduction

# 1. Motivation

- Heavy metal as personal interest
- Distinctive characteristics in visual representation
  - Logos: highly textured patterns, less pictorial
  - Cover arts: without name of band/album; SNS as source
- Visually complex even for human beings → Deep learning



## 2. Functionality and Problem Domain

- Input: 1 real-life photo w/ logo(s) and/or album(s)
- Output: related album/band information (finally returned Metal-archive band profile link)
- Object recognition
- Feature extraction
- Instance matching

## 2. Functionality and Pr

- Object recognition
- Feature extraction
- Instance matching



## 2. Functionality an

- Object recognition
- Feature extraction
- Instance matching



## 2. Functions

- Objects
- Features
- Instances

### MASTODON

Country of origin: United States  
Location: Atlanta, Georgia  
Status: Active  
Formed in: 2000

Genre:

Progressive/Sludge  
Metal (early),  
Progressive Metal/Rock  
(later)

Lyrical themes:

Life, Hardships, Heroic  
tales, Mythical beasts,  
Literature, Spirituality  
Reprise Records

Years active: 2000-present

Current label:

Officially formed on January 13, 2000, a few days after Brann Dailor and Bill Kelliher had met the other members at a High on Fire show.

Founded as a five-piece with Eric Saner on vocals, who left the band just before going on tour. The rest of the band members completed the tour sharing vocal duties, and eventually decided to continue as a four-piece.

Scott Kelly from Neurosis has appeared as ...

{ report  
an error }



READ MORE +

Путь Наверх Часть 1

BAND: <https://www.metal-archives.com/bands/%D0%9A%D0%B8%D0%BF%D0%B5%D0%BB%D0%BE%D0%B2/8649>

Nightblade

Servant To Your Lair

<https://www.metal-archives.com/bands/Nightblade/29763>

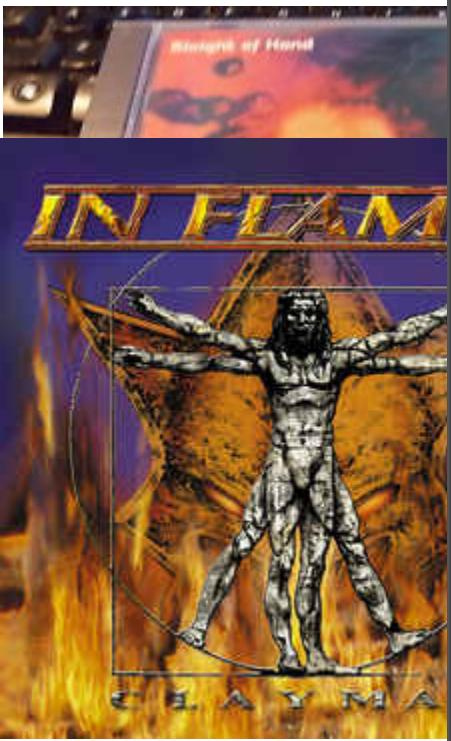
<https://www.metal-archives.com/bands/Nightblade/3540353095>

# Sem 1: Building Pipeline

1



18

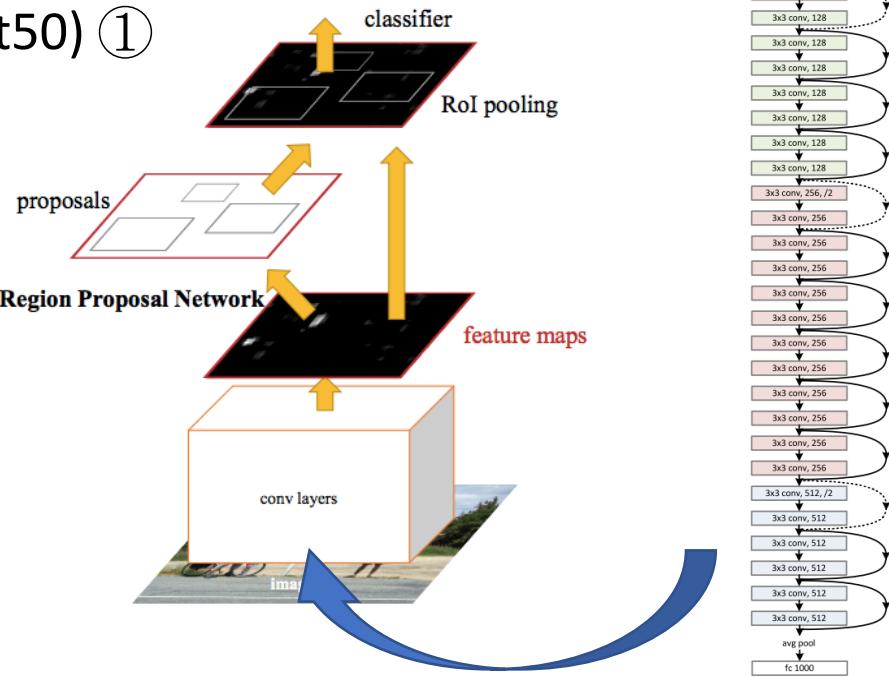


band\_html  
band\_id  
band\_stats  
country  
location  
status  
formed\_in  
genre  
theme  
label  
years\_active  
current\_members  
similar\_artists  
related\_links  
comment  
logo\_url



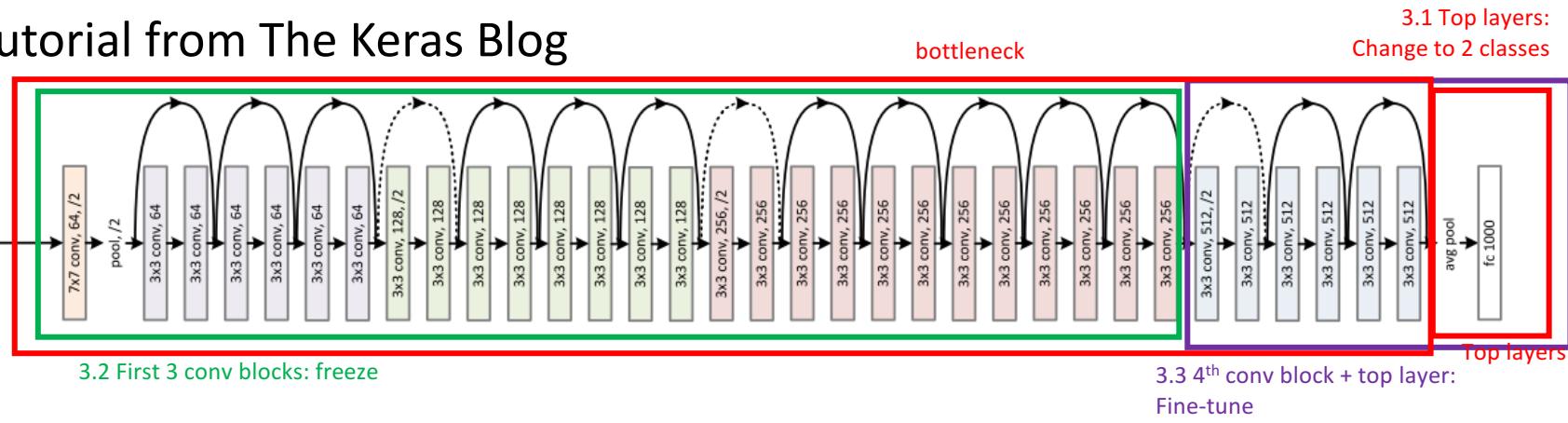
## 2. Object Recognition

- Detection – new FRCNN end-to-end ②
- Classification – pre-train CNN (Resnet50) ①



# Pre-train Resnet50

- Fine-tune top layers of pre-trained ImageNet networks
- Tutorial from The Keras Blog

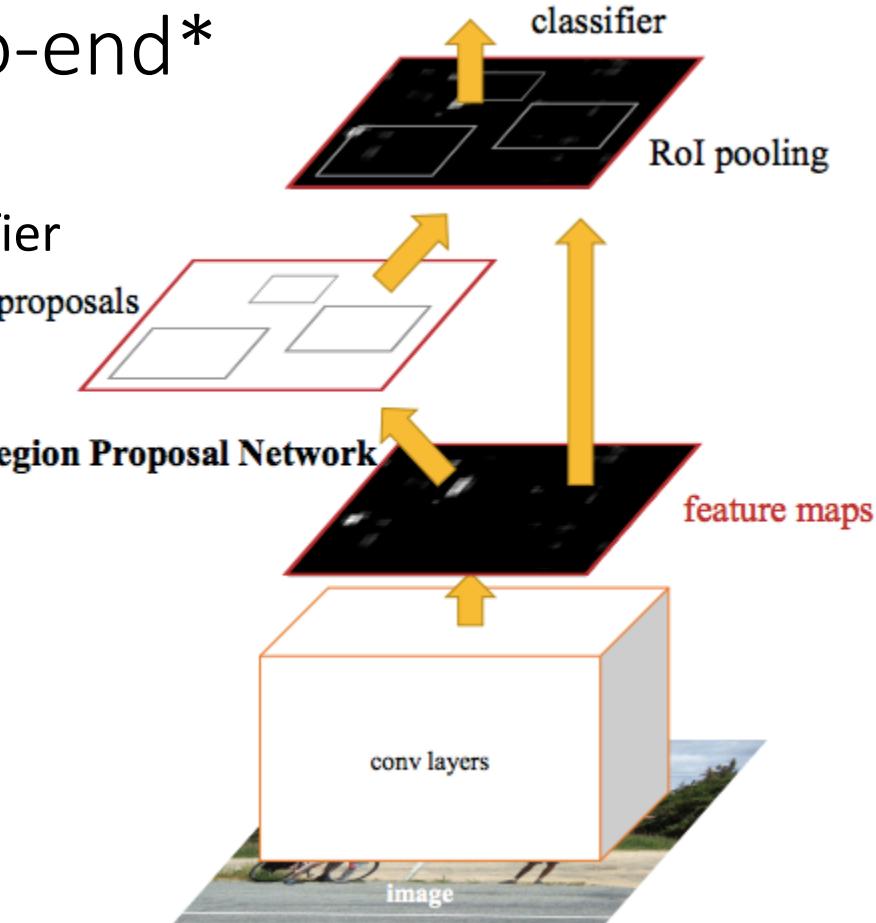
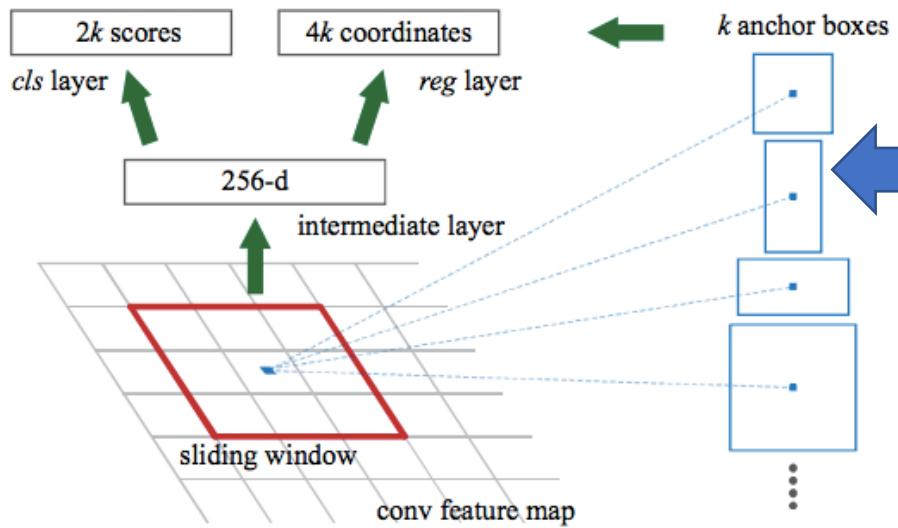


1. Pre-trained model weights → network, run train img through → bottleneck features
2. Train self-def top model (2 classes) → top model weights
3. → fine-tuned model weights

6000 images of each class (12000 in total)  
Training/validation 3:1  
Step 2 and 3 each 100 epochs  
Accuracy > 0.95, loss < 0.01

# Train New FRCNN End-to-end\*

- FRCNN architecture:  
CNN conv layers + RPN + CNN classifier



# Train New FRCNN End-to-end\*

- Train img + bbox info → network → Roi info + obj class



5/135 real241.jpg

● logo

1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18	19	20	21	22	23	24	25	26
1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18	19	20	21	22	23	24	25	26

```
<annotation>
<folder>test1228</folder>
<filename>real241.jpg</filename>
<size>
<width>621</width>
<height>960</height>
</size>
<segmented>0</segmented>
<object>
<name>logo</name>
<bndbox>
<xmin>256</xmin>
<ymin>746</ymin>
<xmax>524</xmax>
<ymax>790</ymax>
</bndbox>
</object>
<object>
<name>logo</name>
<bndbox>
<xmin>99</xmin>
<ymin>737</ymin>
<xmax>239</xmax>
<ymax>798</ymax>
</bndbox>
</object>
```

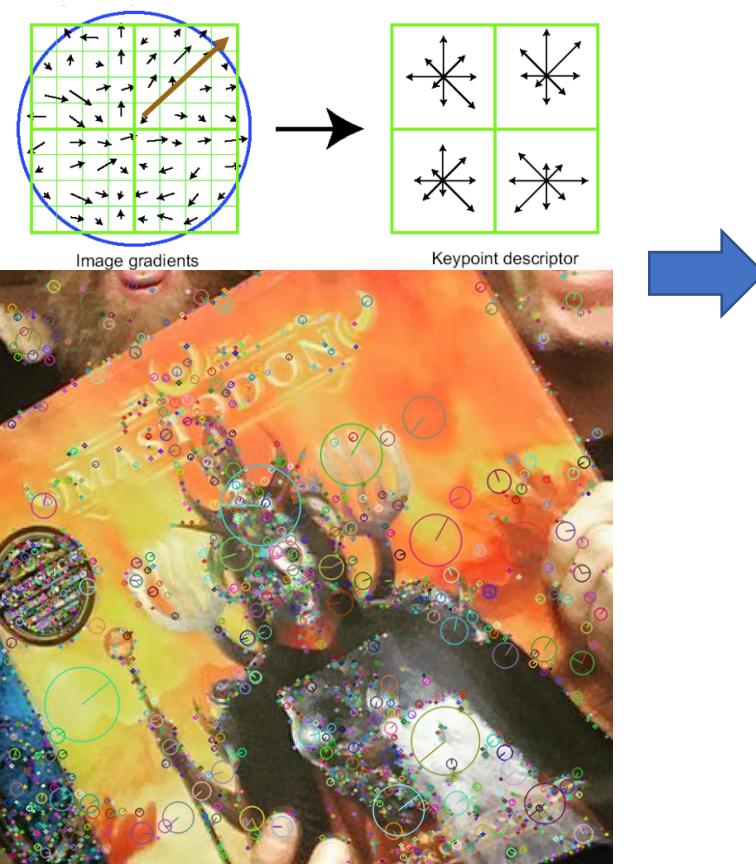
# Train New FRCNN End-to-end\*

- → FRCNN model weights

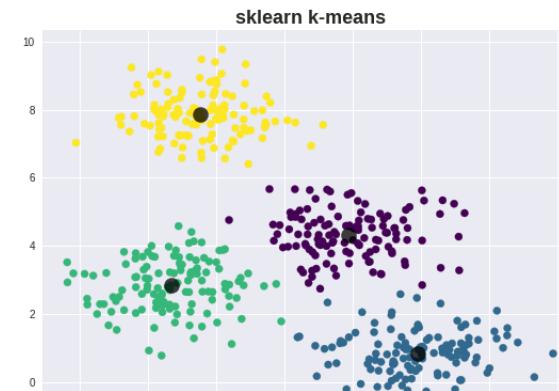
```
Epoch 298/350
999/1000 [=====>.] - ETA: 1s - rpn_cls: 0.2556 - rpn_regr: 0.0404 - detector_cls: 0.0440 - detector_regr: 0.0059Average number of overlapping bounding boxes from RPN = 25.261 for 1000 previous iterations
1000/1000 [=====1 - 1522s - rpn_cls: 0.2557 - rpn_regr: 0.0404 - detector_cls: 0.0440 - detector_regr: 0.0059
Mean number of bounding boxes from RPN overlapping ground truth boxes: 25.239
Classifier accuracy for bounding boxes from RPN: 0.98315625
Loss RPN classifier: 0.281904333636
Loss RPN regression: 0.0405812064885
Loss Detector classifier: 0.0440586159273
Loss Detector regression: 0.00606954138941
Elapsed time: 1522.04713583
Total loss decreased from 0.384222835034 to 0.372613697442, saving weights
Best weights so far saved to ./model6.hdf5. best_loss = 0.372613697442
[[ 5.40990470e-08   3.61807202e-03   9.12488922e-02   7.39599857e-03
  9.68750000e-01]
 [ 1.04356118e-07   4.39091511e-02   5.59981167e-02   2.77232495e-03
  9.68750000e-01]
 [ 2.52584994e-01   1.00336432e-01   3.02801607e-03   6.93191495e-03
  1.00000000e+00]
 ...,
 [ 5.90860658e-08   2.74258386e-03   1.20253884e-03   4.21987521e-03
  1.00000000e+00]
 [ 6.19410301e-08   1.87421683e-02   1.69737056e-01   3.37358974e-02
  9.68750000e-01]
 [ 5.59893074e-08   6.67789802e-02   1.09739833e-01   2.01013666e-02
  9.68750000e-01]]
```

1312 train img  
500 epochs in total

### 3. Feature Extraction\*



- Keypoint level: SIFT descriptors,  $10^2 \sim 10^3$  per img
- k-means cluster SIFT from all imgs
- Img level: VLAD descriptor, 1 per img
- → ground truth VLAD descriptors



For each centroid  $c_i$ , a sub vector  $v_i$  is obtained by accumulating the residual vector between the centroid and local features quantized to this centroid:

$$v_i = \sum_{x:q(x)=c_i} x - c_i \quad (2)$$

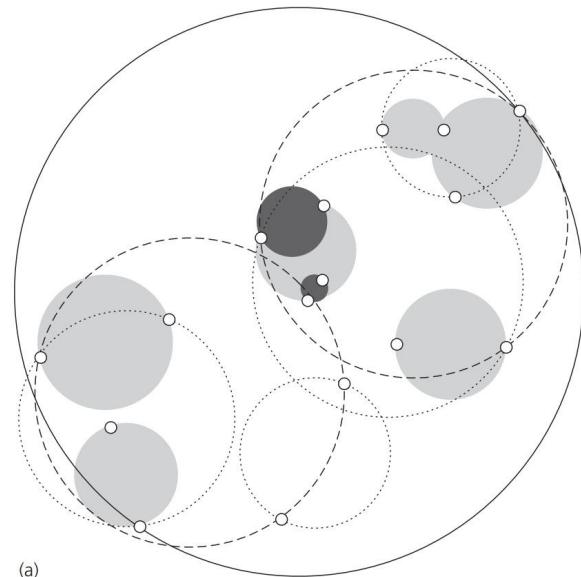
# Initial Experiments on Merged Features

- Purpose for initial test is only for illustration and finding potential better-performance feature combinations.

No.	Merged Features	Successful Rate %
i	VLAD	100
ii	VLAD + Object-level FRCNN features	22
iii	Object-level FRCNN features	22
iv	FRCNN column features	22
v	Object-level + FRCNN column features	22
vi	VLAD + Object-level features of Pre-trained CNN	56
vii	VLAD + Column features of Pre-trained CNN	67
viii	Object-level features of Pre-trained CNN	56
ix	Column features of Pre-trained CNN	67
x	Object-level + Column features of Pre-trained CNN	33
xi	VLAD + Color Histogram	0
xii	Color Histogram	0
xiii	VLAD with PCA (6400 components)	<30 (didn't keep record)
xiv	VLAD with PCA (6400 components) + Object-level features of Pre-trained CNN	<30 (didn't keep record)
xv	VLAD with PCA (12800 components)	<30 (didn't keep record)

## 4. Instance Matching

- VLAD descriptors → construct ball trees
- “comprehensive search for 3D space”
- Search for 10 most similar imgs, return IDs



# 5. Information Retrieval

- Info from 2 online databases → 2 tables
- Objective: select URLs where ID = ID retrieved in last step

Album

hot4sql		ma4sql	Data	Schema	SQL
Tables	Album	Band	ImgID		
hot4sql 555 Rows	1 Phantom Anthem	August Burns Red	R-10960429-1...7241009-9781		
ma4sql 116,887 Rows	2 Leviathan	Alestorm	R-11043047-1...8784572-9240		
	3 Axe Crazy	Jaguar	R-1867902-1289061630		
	4 Sweet Death And Ecstasy	Midnight	R-11037797-1508693751-5740		
	5 Metallurgy2 -...ns To Be Fearful	Various	R-1830191-1246306636		
	6 Today We Become The Enemy	The 244GL	R-2050557-1260915676		

Logo/bandinfo

hot4sql		ma4sql	Data	Schema	SQL
Tables	ID	Name	URL		
hot4sql 555 Rows	1 3540410023	A B I S M O	https://www.m...O/3540410023		
ma4sql 116,887 Rows	2 3540352307	A Balance of Power	https://www.m...r/3540352307		
	3 108513	A Band of Orcs	https://www.m...of_Orcs/108513		
	4 3540412046	A Baptism by Fire	https://www.m...e/3540412046		
	5 18563	A Blind Prophecy	https://www.m...rophecy/18563		
	6 111183	A Bloody Epitaph	https://www.m...Epitaph/111183		
	7 103063	A Body Falls	https://www.m...y_Falls/103063		
	8 3540371582	A Breach of Silence	https://www.m...e/3540371582		
	9 3540333609	A Break in the Storm	https://www.m.../3540333609		
	10 121141	A Breath Before Surfacing	https://www.m...urfacing/121141		

10,000 records in 0.034 seconds

< 1 of 12 >

# In conclusion, query has 5 steps:

Target	Content	Input	Output/Mid product
1 Query img	Obj rec in FRCNN	Query img	Cropped Rols (O)
1 Crop	VLAD extraction*	Crop	SIFT(VLAD) descriptor of crop (M)
1 Crop	Instance matching	SIFT(VLAD) of crop	Matched img IDs (M)
1 Crop	Info retrieval	Matched img IDs	(temp) Possible band URLs (M)
1 Query img	Info retrieval	URLs of each crop	All possible band URLs (O)

# Initial End-to-end Test

- Precision of object recognition and overall accuracy were calculated wrong in the mid-point presentation, actual values are lower.

Item	Remark	Number
Album ground truth image set		500
Logo ground truth image set		500
Test images		155
All crops	Crops = Rols written to file	426
Correct crops	Crops containing album or logo	316
Matched crops	Correct band information retrieved	212
Recall of object recognition	Recall = correct Rols/correct Rols+Rols haven't been recognized (multiple crops on one object is counted once)	159/192 = 83.68
Precision of object recognition	Precision = correct crops/(correct+wrong) crops	170/426 = 39.91
Overall accuracy		33.40



position  
better



# Sem 2: Optimization

# 1. Network Optimization

- Train a new FRCNN -> Finetune pretrained ImageNet FRCNN
- Add cross validation
- Add data augmentation
- Human eval -> comprehensive mAP calculation
- Generalization ability: balance between underfitting/overfitting

# 1. Network Optimization

- Train a new FRCNN -> Finetune pretrained ImageNet FRCNN
  - FRCNN = Resnet50 conv layers + RPN + Resnet50 classifier
  - Sem 1 used self-trained weights for conv, RPN/classifier trained from zero
  - Sem 2 used self-trained weights for conv/classifier, pretrained weights for RPN, then finetuned
- Add cross validation
- Add data augmentation
- Human eval -> comprehensive mAP calculation
- Generalization ability: balance between underfitting/overfitting

# 1. Network Optimization

- Train a new FRCNN -> Finetune pretrained ImageNet FRCNN
- Add cross validation
  - Matriculated with the source code, found out the cross validation section was dysfunctional, modified
- Add data augmentation
- Human eval -> comprehensive mAP calculation
- Generalization ability: balance between underfitting/overfitting
  - Diagnostic stats/best weights saved from mAP calculation whenever loss decreases for plotting training loss vs. test loss

# 1. Network Optimization

- Train a new FRCNN -> Finetune pretrained ImageNet FRCNN
- Add cross validation
- Add data augmentation
  - Introduced noise
- Human eval -> comprehensive mAP calculation
- Generalization ability: balance between underfitting/overfitting
  - Diagnostic stats/best weights saved from mAP calculation whenever loss decreases for plotting training loss vs. test loss

# 1. Network Optimization

- Train a new FRCNN -> Finetune pretrained ImageNet FRCNN
- Add cross validation
- Add data augmentation
- Human eval -> comprehensive mAP calculation
  - Sem 1: used OpenCV imshow to check the Rols, manually counted, bound to error
  - Sem 2: wrote script for auto calculation
- Generalization ability: balance between underfitting/overfitting

# 1. Network Optimization

- Train a new FRCNN -> Finetune pretrained ImageNet FRCNN
- Add cross validation
- Add data augmentation
- Human eval -> comprehensive mAP calculation
- Generalization ability: balance between underfitting/overfitting
  - Diagnostic stats/best weights saved from mAP calculation whenever loss decreases for plotting training loss vs. test loss

170/426 =  
39.91  
33.40

# Trials after Adjusting Network Configuration

Trials	AP album	AP logo	mAP	loss	Rank
finetune + cross val + final mAP	0.875773854135	0.685576560585	0.78067520736	1.34752395183	2
finetune + aug + final mAP	0.75902475123	0.696653130193	0.727838940712	2.66465598715	6
finetune + cross val + aug + final mAP	0.871920340588	0.667815494874	0.769867917731	1.42401295861	4
finetune + cross val + per optimal mAP + gen/plt (best)	0.893460460836	0.727167063276	0.810313762056	1.2770022666	1
finetune + cross val + per optimal mAP + gen/plt (2 <sup>nd</sup> best)	0.814064204169	0.733884851467	0.773974527818	1.37275927728	3
finetune + cross val + aug + per optimal mAP + gen/plt (best)	0.846656380608	0.645819069931	0.74623772527	1.55319615773	5
finetune + cross val + aug + per optimal mAP + gen/plt (2 <sup>nd</sup> best)	0.727294211281	0.688764544586	0.708029377933	2.89686328396	7

# (Opt) Hyperparameter Optimization

- Python lib: hyperopt, idea: random search
- Hyperspace, Fmin
- Implementation: training 20 epochs to save time, output config in JSON, retraining
- Choice of hyperparameters: lr, decay, kernel size, anchor box scale
- Problem and insight
- Conclusion: after retraining, best model with **0.81** accuracy

# Demo: Re-eval of FRCNN



## 2. Problem on Feature Fusion

- Inappropriate input img size (224,224) -> (600,?)/(?,600) keep the ratio

<b>Step 1, branch 1</b>	SIFT	a. Describe SIFT features on all training images b. Kmeans clustering on all SIFT descriptors from all training images to build visual dictionary c. Calculate VLAD vector for each image based on visual dictionary
<b>Step 1, branch 2</b>	CNN	d. Extract conv5 feature as object level feature e. Extract conv4 feature as column feature <ul style="list-style-type: none"><li>• Misunderstanding of column feature/FC feature</li></ul>
<b>Step 2</b>	Concatenation of CNN feature and VLAD descriptor	
<b>Step 3</b>	Index ball tree for searching	

# Trials after Changing Input Img Size

	Album		Logo	
	Dimension	Accuracy	Dimension	Accuracy
SIFT(VLAD)	65536->65536	0.8877551	65536->65536	0.39215686
CNN object level feature (original weights)	2048->8192	0.41326531	2048->28672	0.33529412
CNN object level feature (model 1)	2048->8192	Fail	2048->28672	Fail
CNN object level feature (model 2)	2048->8192	Fail	2048->28672	Fail
SIFT(VLAD) + CNN object level feature (original weights)	67584->73728	0.41326531	67584->94208	0.33529412
SIFT(VLAD) + CNN object level feature (model 1)	67584->73728	Fail	67584->94208	No trial
SIFT(VLAD) + CNN object level feature (model 2)	67584->73728	Fail	67584->94208	No trial

- Higher overall than Sem 1.
- SIFT on album could be considered as a successful case.

# Trials with FC Features

	Accuracy
FC	0.18823529
SIFT(VLAD) + FC	0.12941176
SIFT(VLAD) with PCA + FC	0.22745098

- Sub-ideal result

# New Mechanism: VLAD on CNN conv5

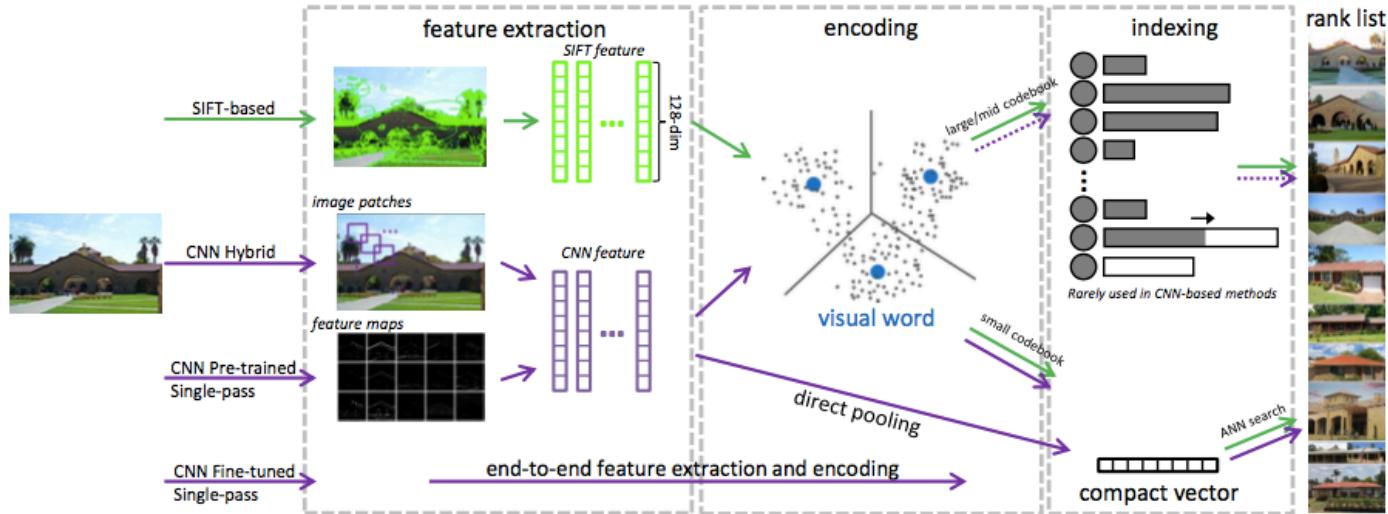


Fig. 2: A general pipeline of SIFT- and CNN-based retrieval models. Features are computed from hand-crafted detectors for SIFT, and densely applied filters or image patches for CNN. In both methods, under small codebooks, encoding/pooling is employed to produce compact vectors. In SIFT-based methods, the inverted index is necessary under large/medium-sized codebooks. The CNN features can also be computed in an end-to-end way using fine-tuned CNN models.

# Trials with VLAD on CNN conv5

	Accuracy
CONV5(VLAD) on hard images	0.33529412
SIFT(VLAD) + CONV5(VLAD) on hard images	0.44313725
SIFT(VLAD) + CONV5(VLAD) on FRCNN Rols	0.64705882

- At least... it's better.

# (Opt) New Mechanism: Canonical Correlation Analysis

- Sub-ideal results.

# Demo: Re-eval of Logo Matching

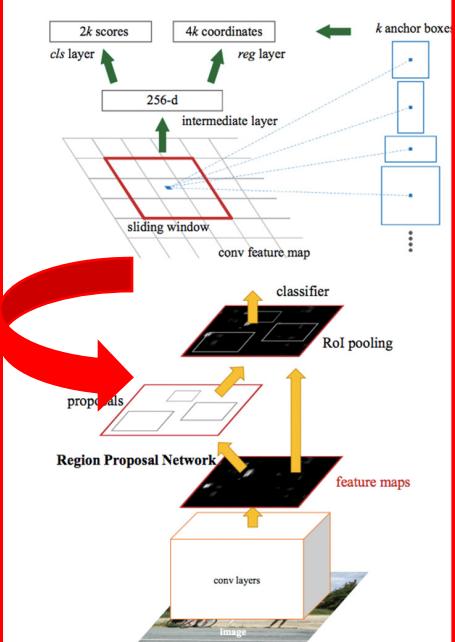
- Conclusion:
- SIFT(VLAD) + CNN conv5(VLAD) for logos, balancing the time consumption and accuracy.
- SIFT(VLAD) for albums.
- Best mAP  $\approx 0.76$ .
- **End-to-end accuracy  $\approx 0.65$** , reasonable considering FRCNN and CNN each have an accuracy around 0.80.

# Finalized Server Pipeline

## 1. Object Recognition – FRCNN

Detection – RPN

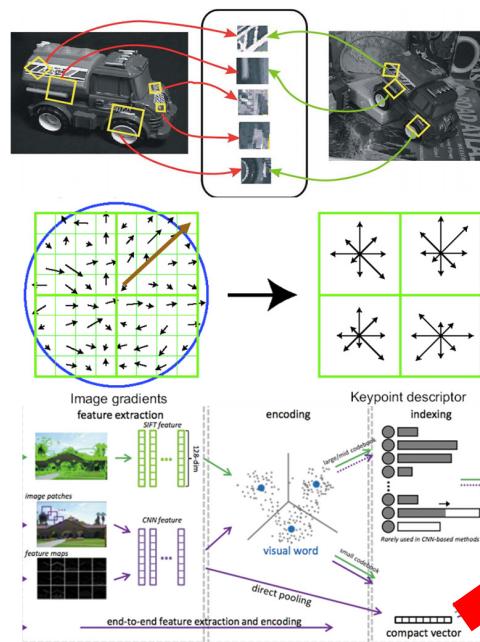
Classification – CNN classifier



## 2. Feature Extraction

SIFT on albums

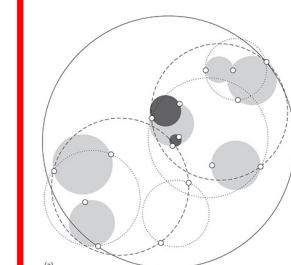
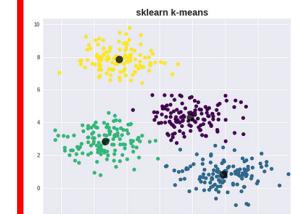
SIFT + conv5 on logos



## 3. Instance Matching

Words – Kmeans

Indexing – Ball Tree



## 4. Information Retrieval

Database

Two screenshots of a database interface showing tables "hot4sql" and "ma4sql".

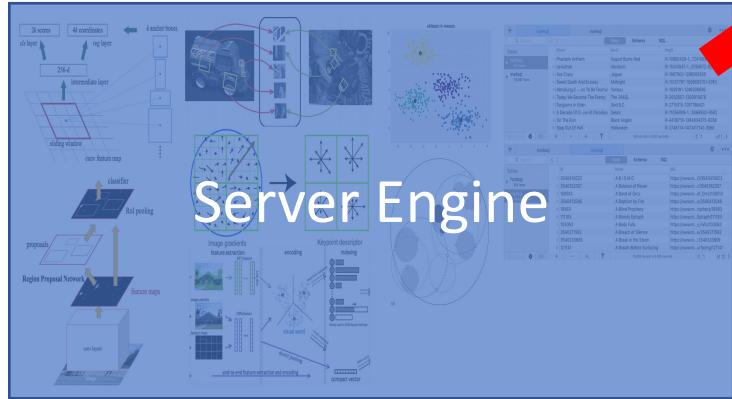
**hot4sql** table data:

ID	Name	URL
1	August Burns Red	R-10960429-1...7241009-9781
2	Leviathan	R-11043047-1...8784572-9240
3	Axe Crazy	R-167702-1289061630
4	Sweet Death And Ecstasy	Jaguar
5	Metallurgy2 -_ns To Be Fearful	Midnight
6	Today We Become The Enemy	R-11037797-1508693751-5740
7	Eargasms In Eden	Various
8	A Decade Of D...ive At Paradise	The 244GL
9	On The Run	R-1830191-1246306536
10	Step Out Of Hell	R-2050557-1760915676

**ma4sql** table data:

ID	Name	URL
1	A B I S M O	https://www.m.../0/3540410023
2	A Balance of Power	https://www.m.../r/3540352307
3	A Band of Orcs	https://www.m.../_of/Orcs/108513
4	A Baptism by Fire	https://www.m.../e/3540412046
5	A Blind Prophecy	https://www.m.../rophecy/18563
6	A Bloody Epitaph	https://www.m.../_Epitaph/11183
7	A Body Falls	https://www.m.../Falls/103063
8	A Breach of Silence	https://www.m.../121141
9	A Break in the Storm	https://www.m.../3540333609
10	A Breath Before Surfacing	https://www.m.../surfacing/121141

# Sem 2: Mobile Development (Opt) MVC Design



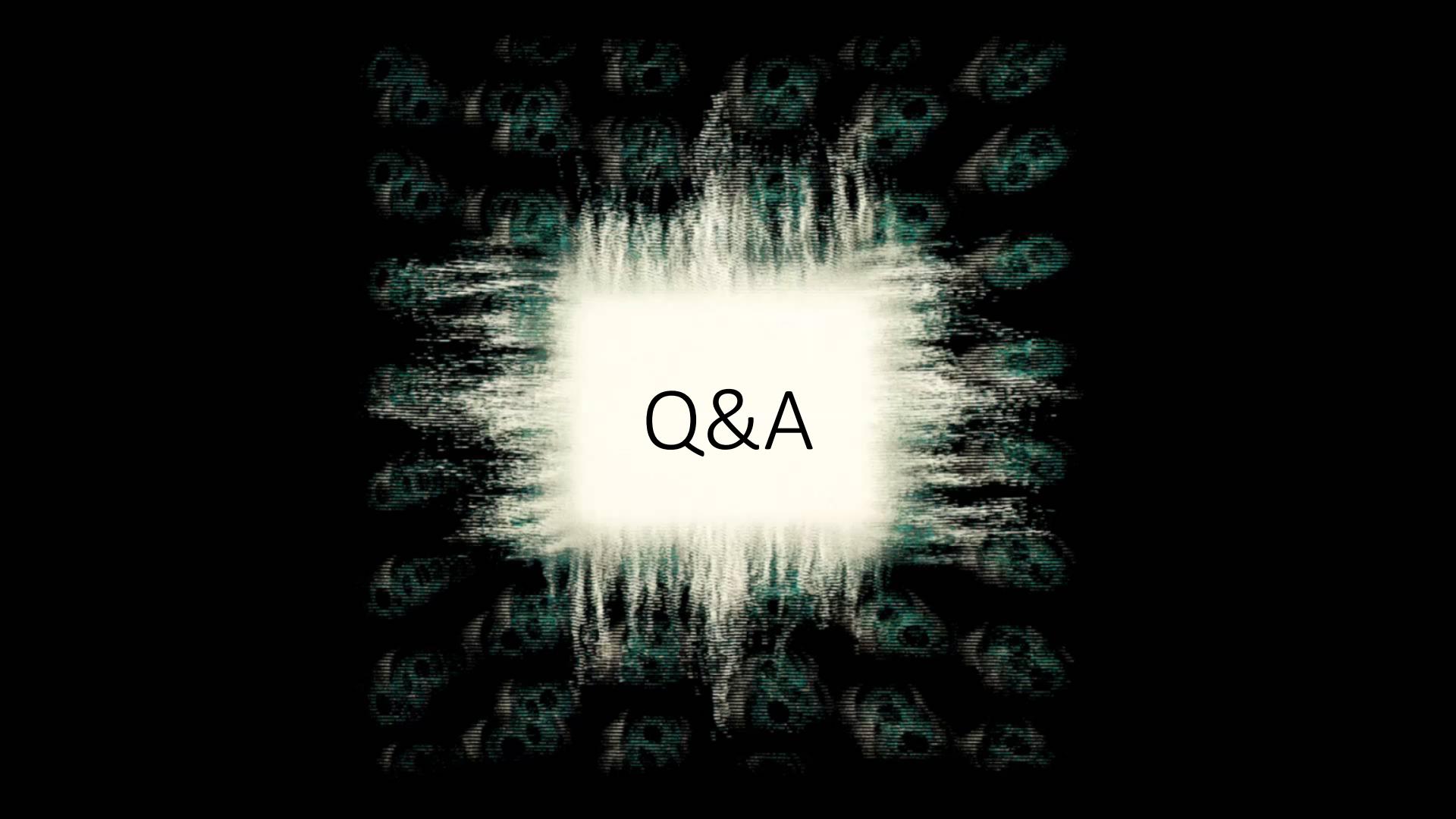
Retrieved Information

JSON

Network



Demo: the APP



Q&A