# About

- Established in 2017
- Apply AI technology to games
- Research Interests: CV, NLP, RL and Speech Signal Processing
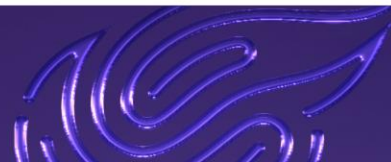
# Outline

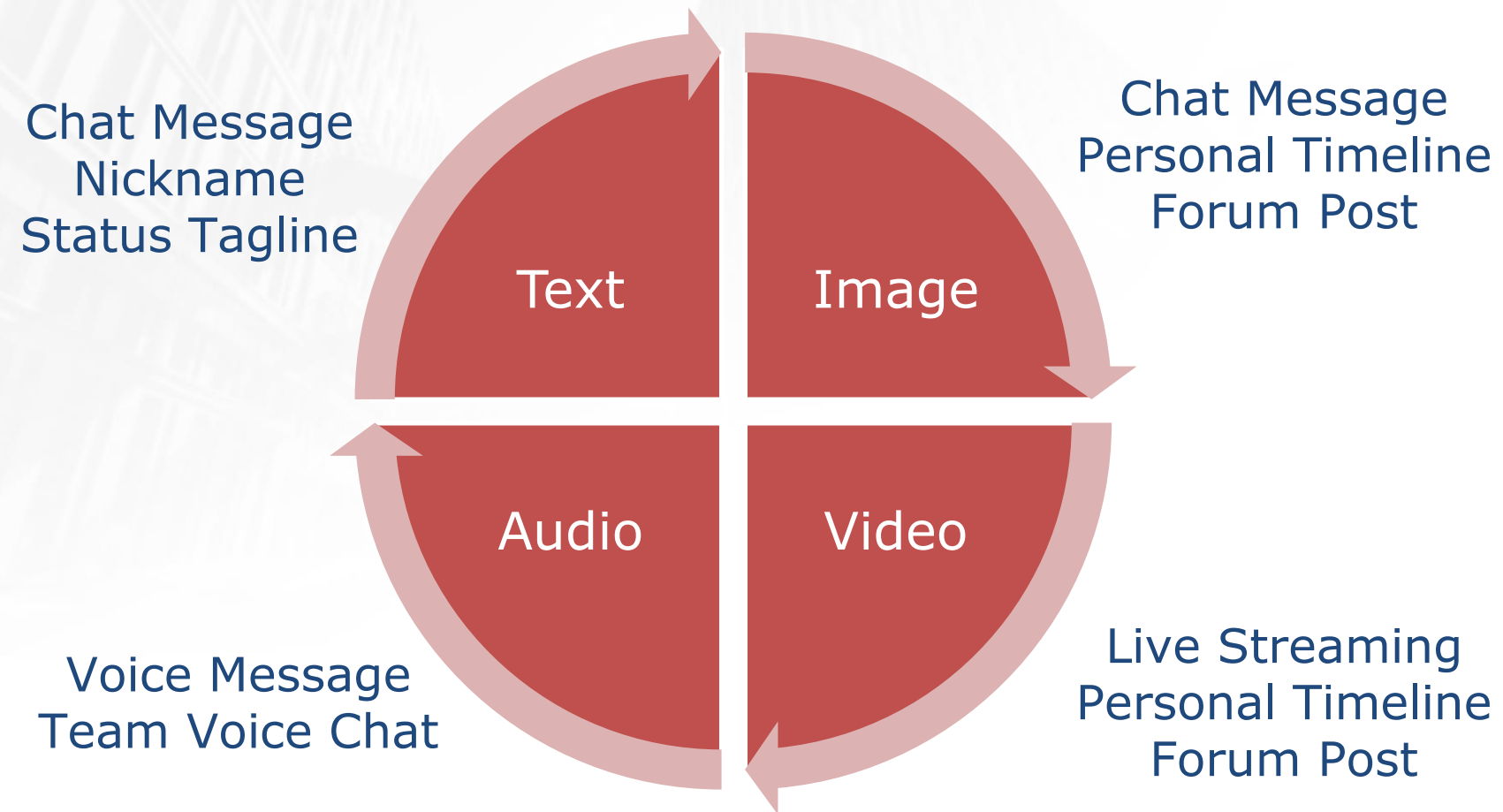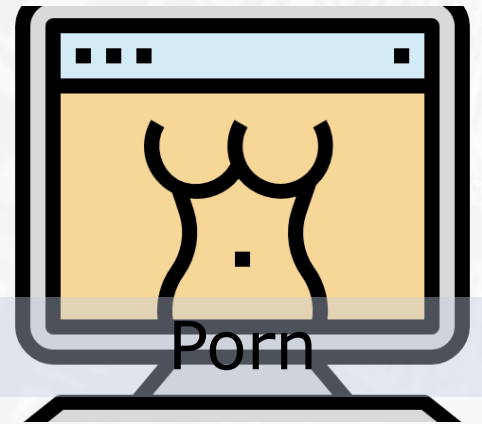GDC × NetEase Games *Passion of gamers*  Netease Games AI Lab

**March 21-25, 2022**
San Francisco, CA

# Introduction



Chat Message
Nickname
Status Tagline

**Text**

Chat Message
Personal Timeline
Forum Post

**Image**

Voice Message
Team Voice Chat

**Audio**

**Video**

Live Streaming
Personal Timeline
Forum Post

**Violence**

**Porn**

**Offensive Content**

**Spam**

**Abuse**

# Introduction



Porn | Sexy | Normal

**Porn**



**Spam**



Weapon | Blood | Terrorist

**Violence**



**Abuse**

# Introduction

## Manual moderation

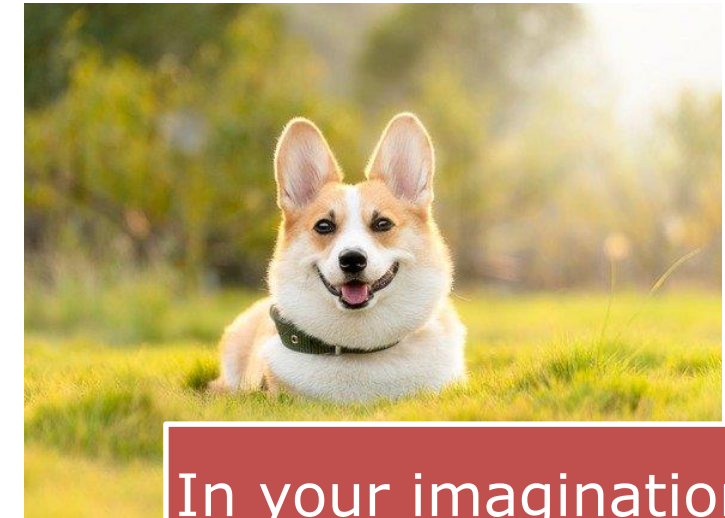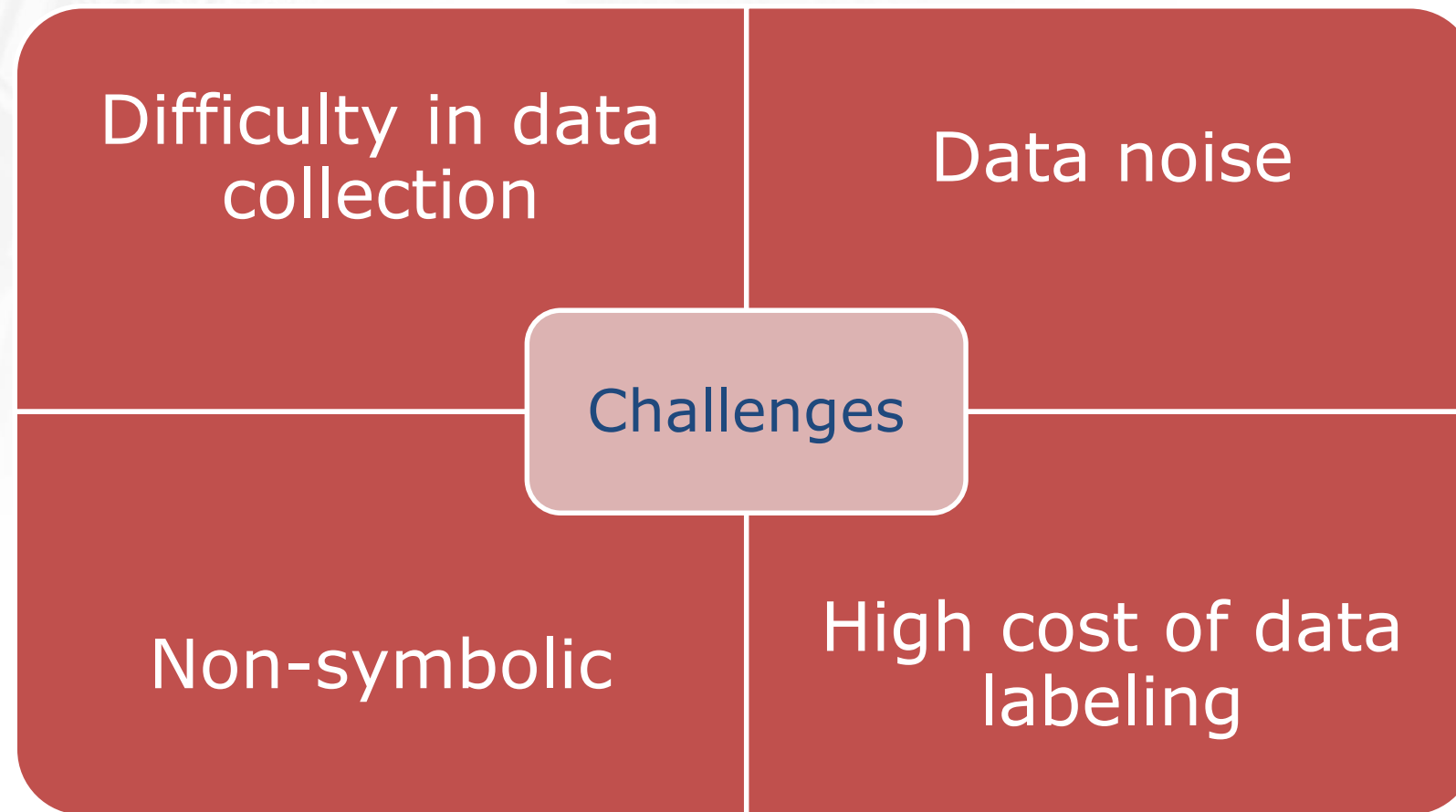| Low efficiency | Not in real-time | Unstable | High cost |
|---|---|---|---|
| Slow manual processing | No 24/7 or immediate response | Inconsistent standards | Requires many moderators |

# Introduction

## Challenges during development


In your imagination


The actual situation

**Non-symbolic image**

| | |
|---|---|
| Difficulty in data collection | Data noise |
| Non-symbolic | High cost of data labeling |

Challenges
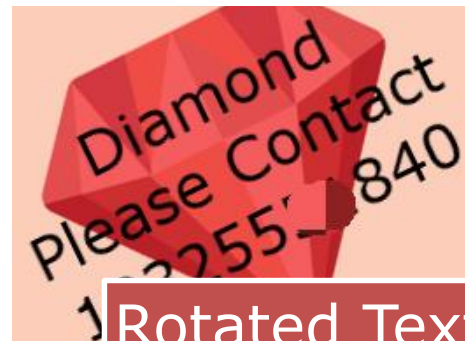
# Introduction

## Challenges during maintenance

New types emerge over time

"Attackers" create new variants

Rotated Text

Curved Text

Handwriting

New variants of image spam
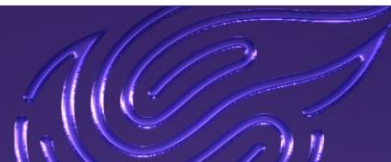
# Outline

- Introduction
- **Image Sub-system**
- Audio Sub-system
- Text Sub-system
- Application

# Define Categories

Create categories for common scenarios and ignore outliers
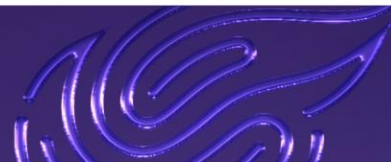
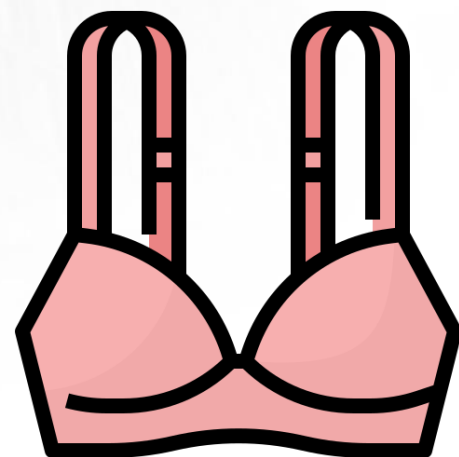Define the scope of each category for precise annotation

Use fine-grained categories to cover more online cases

# Define Categories
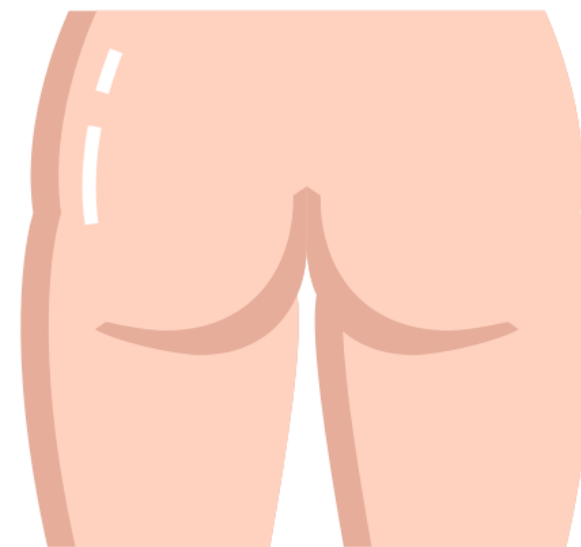
Cleavage

Bra

Underwear

Hip

Fine-grained categories for live streaming

# Data Collection



Web Crawler

Content-based Similarity Search

Internal data accumulation

Data provided by product teams

Update model with online data

Important

# Data Cleaning

| Duplicated Data Removal | Simple Filtering | Data Clustering | Crowdsource Labeling | Deep Cleaning |
|---|---|---|---|---|
| • MD5<br>• pHash/aHash… | • Traditional Features<br>• Data Labeling & SVM/LR Classifiers | • FC Features Clustering<br>• Select dense and compact classes | • Vote by different tools/models<br>• Manual labeling | • K-fold Cross Validation<br>• Inconsistent Predictions Labeling |

# Data Cleaning

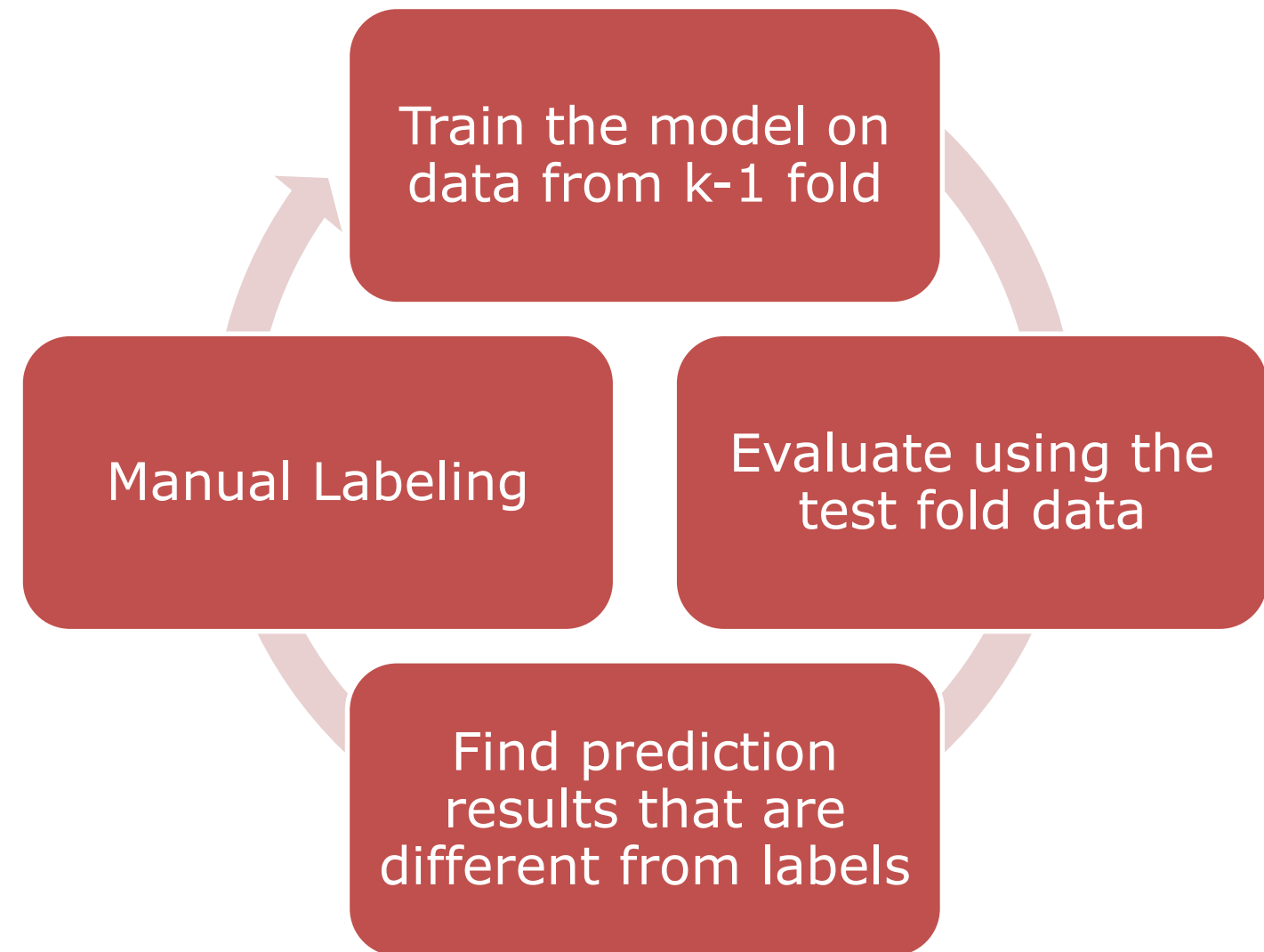| Test | Train | Train | Train | Train |
|------|-------|-------|-------|-------|

| Train | Test | Train | Train | Train |
|-------|------|-------|-------|-------|

| Train | Train | Test | Train | Train |
|-------|-------|------|-------|-------|

| Train | Train | Train | Test | Train |
|-------|-------|-------|------|-------|

| Train | Train | Train | Train | Test |
|-------|-------|-------|-------|------|

Train the model on data from k-1 fold

Evaluate using the test fold data

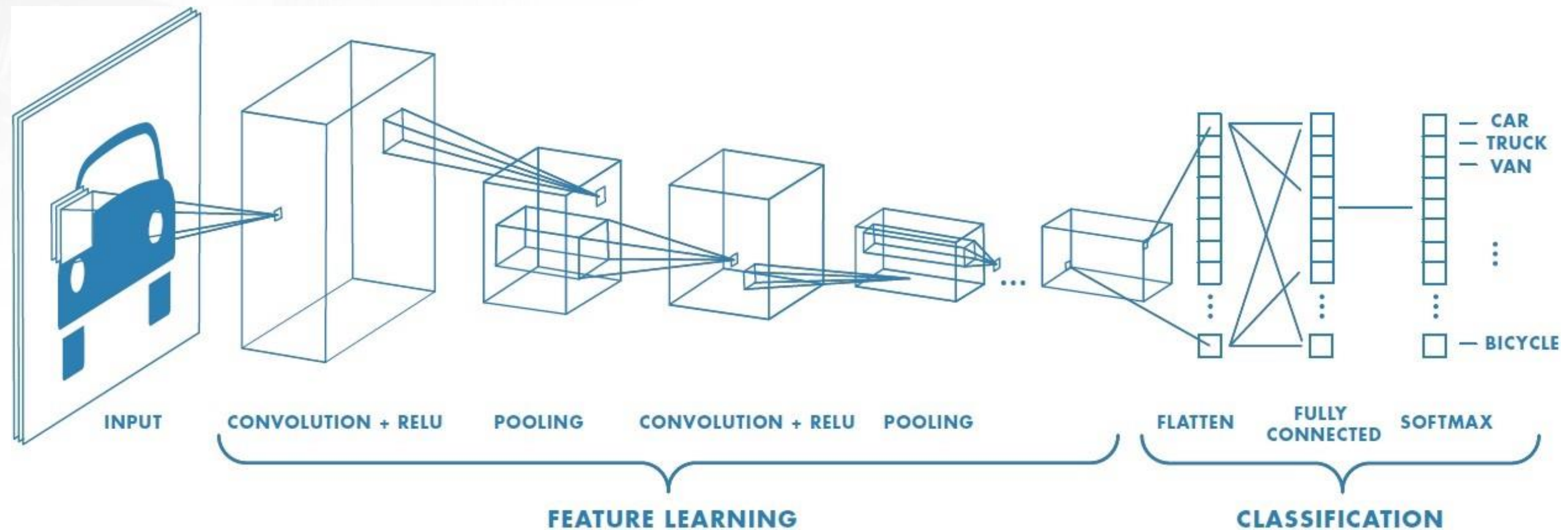Find prediction results that are different from labels

Manual Labeling

5-fold cross validation to get deep cleaned data

# Methods

- Image Classification Networks: ResNet、MobileNet v2/v3



[1] Image source: https://penseeartificielle.fr/mobilenet-reconnaissance-images-temps-reel-embarque/
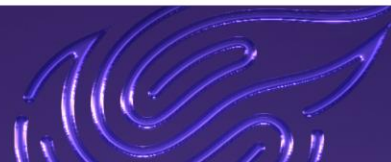
# Methods

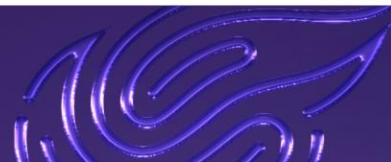| Data Augmentation | Model Overfitting | Imbalanced Classes |
|---|---|---|
| • Auto Augment<br>• Cutout | • Dropout after GAP Layer<br>• Label Smoothing | • Over-sampling<br>• Weighted Loss |

# Model Optimization

Class Activation Map (CAM)

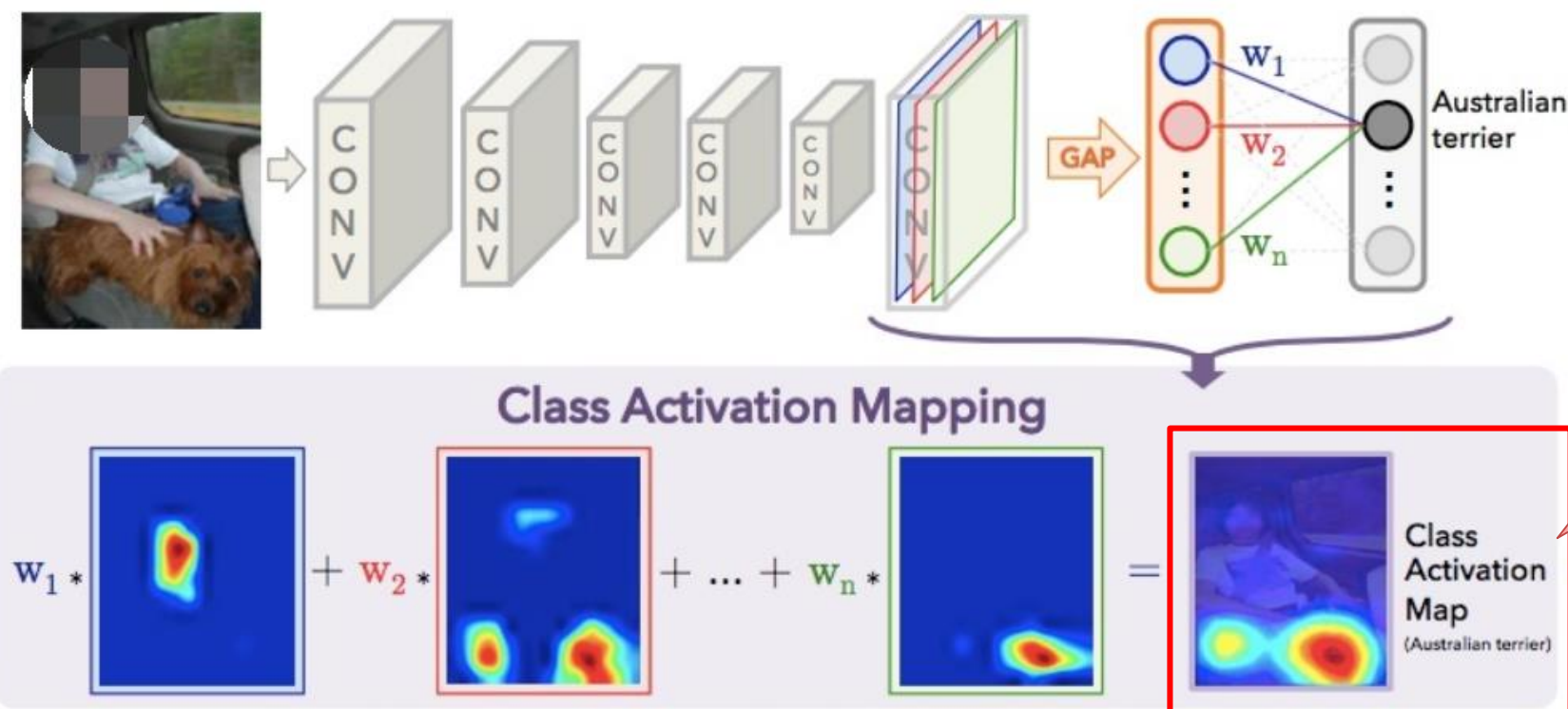Big Input, Small Network

Attention Mechanism

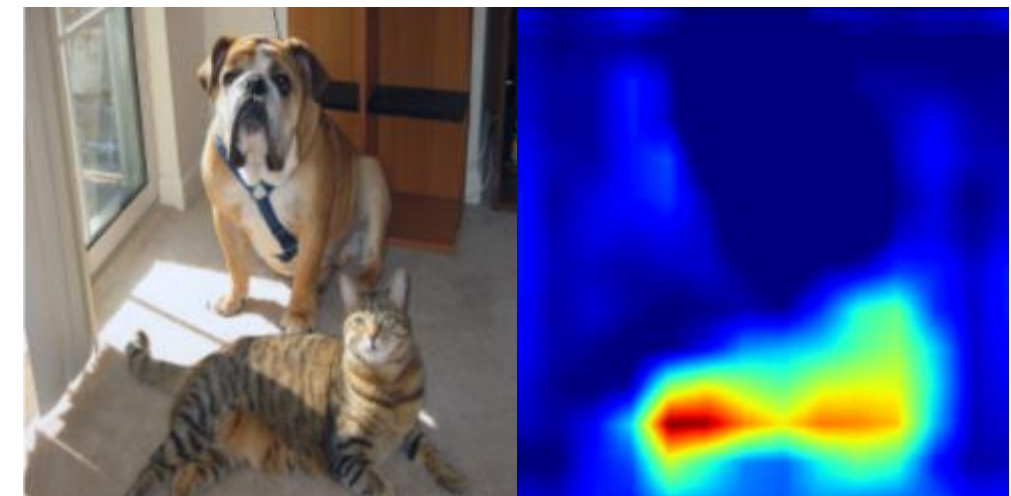Object Detection Networks with FPN: YOLO、RetinaNet

# Model Optimization



Class Activation Mapping

$w_1 * \quad + w_2 * \quad + \dots + w_n * \quad = $

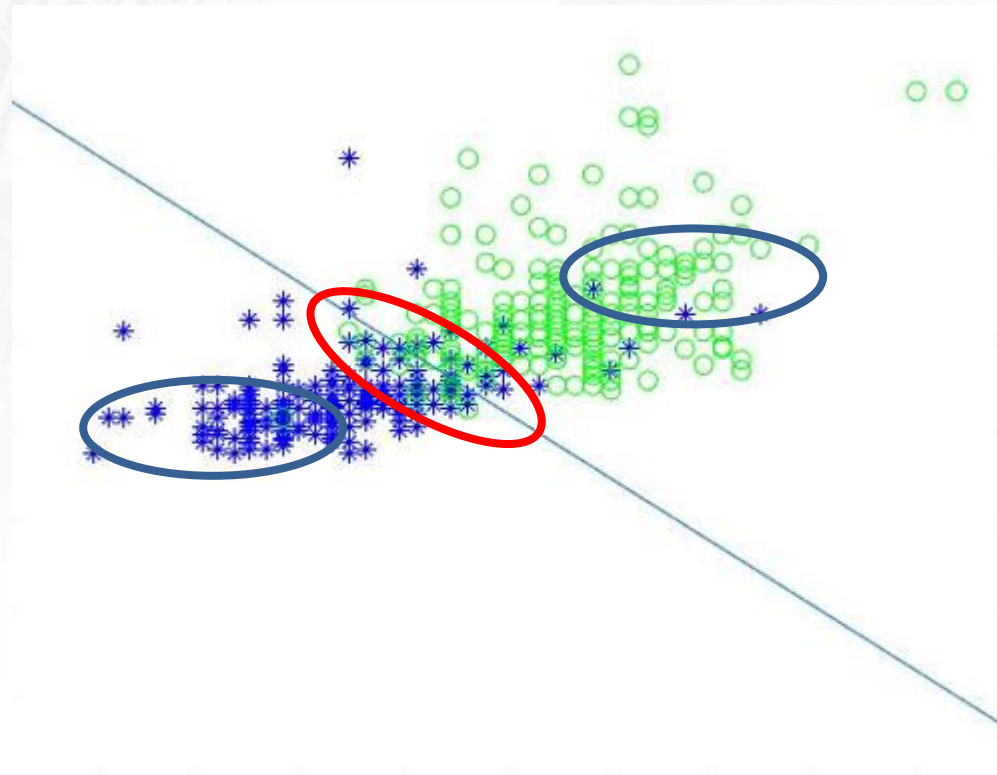Class Activation Map (Australian terrier)

Result

Why the prediction result is **not a dog**?
Ah, the reason is **the influence of cat**.
So we need **collecting more images containing cat**

Use Class Activation Map(CAM) to Improve Data

[1] Zhou B, Khosla A, Lapedriza A, et al. Learning deep features for discriminative localization.

# Model Update



Non-ambiguous data:
Used to increase data diversity
and model stability

Boundary data:
Used to update the data boundary
and increase model robustness

Pulling Additional Online Data to Update the Models
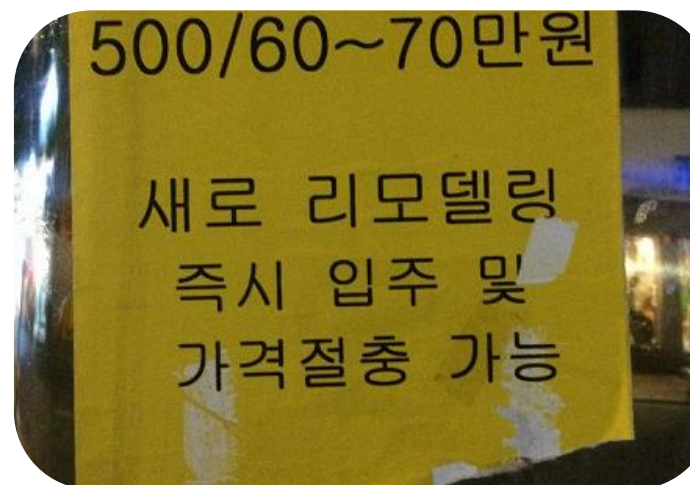
# Scene Text Detection

- Self-labeling datasets and public datasets
- A dataset containing 63,000+ images was constructed
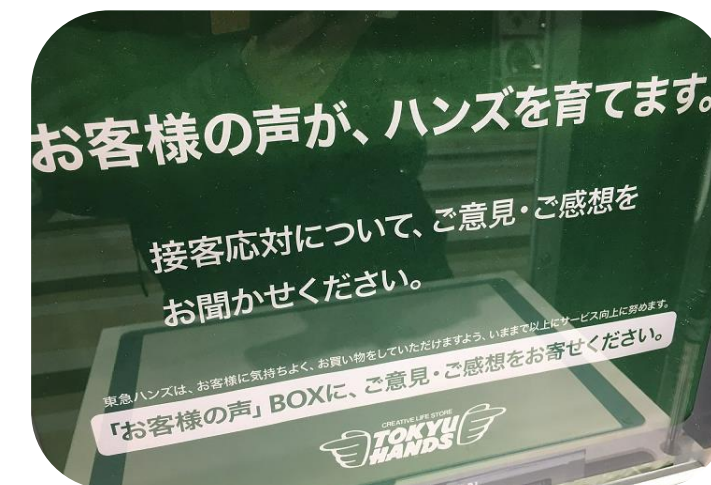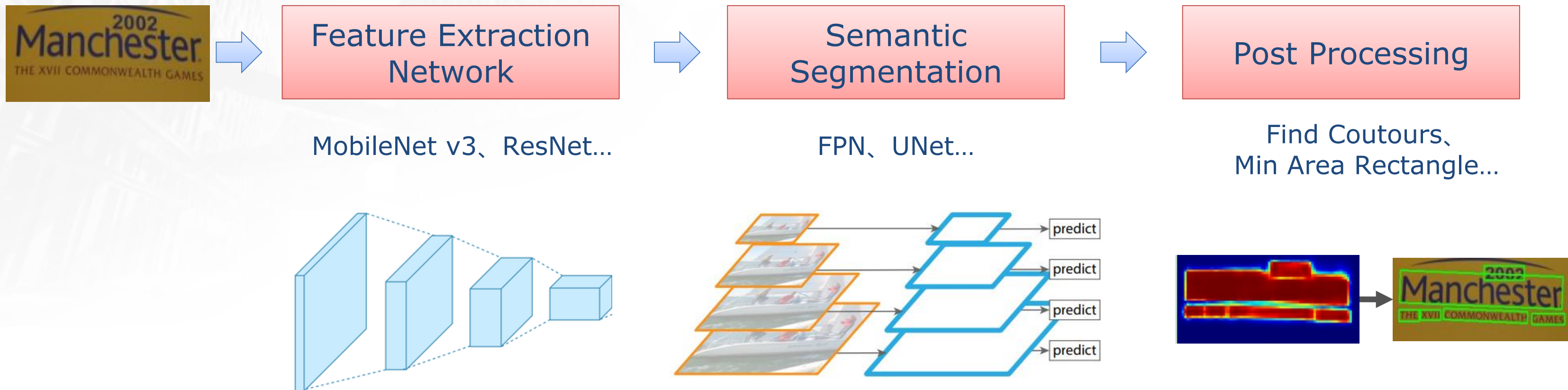


Chinese



English



Korean



Japanese

# Scene Text Detection



Feature Extraction Network

MobileNet v3、ResNet...

Semantic Segmentation

FPN、UNet...

Post Processing

Find Coutours、
Min Area Rectangle...
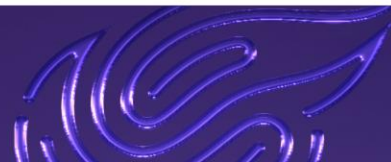
predict
predict
predict
predict

# Scene Text Recognition

- Synthesis dataset
- Real-world dataset annotated by crowdsource labeling of several tools/models

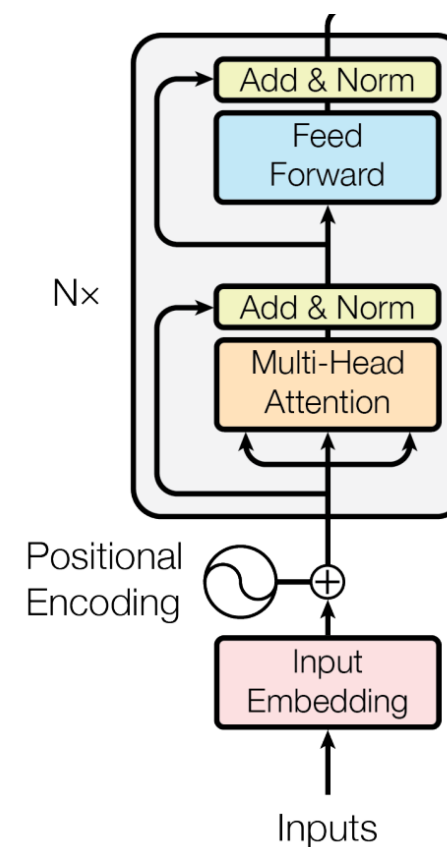douses nonimperiously fluorid macrophotograph masturbatic
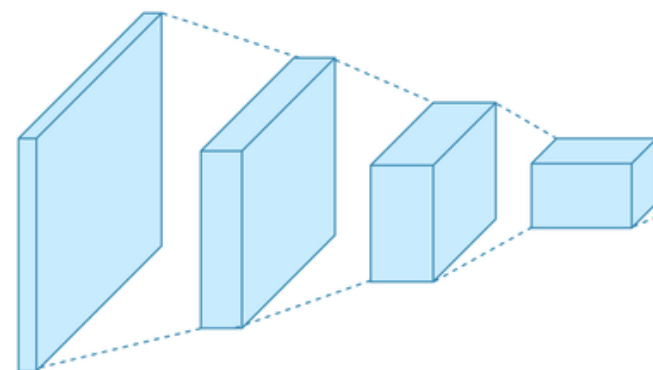
EMMALINE REDEVABLE LECKKILL RETRORENAL CURTSY'S

expeditionary synetic palatograph mementoes cordwainer

marlins Lemonias bookbinder Spondylus jelled

# Scene Text Recognition

| Pre Processing Network | Image Feature Extraction Network | Sequential Feature Extraction Network | Decoder |
|---|---|---|---|
| STN、TPS... | ResNet、RepVGG... | Transformer Encoder、LSTM... | CTC、Attention... |

# Outline

- Introduction
- Image Sub-system
- **Audio Sub-system**
- Text Sub-system
- Application

# Speech Recognition

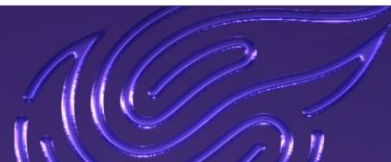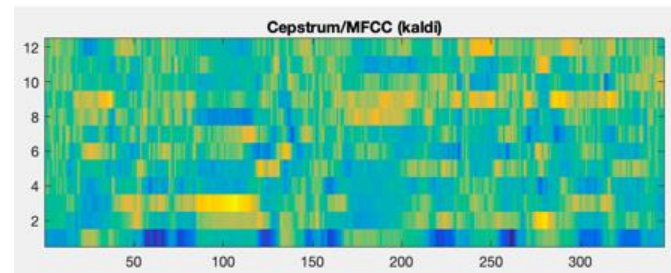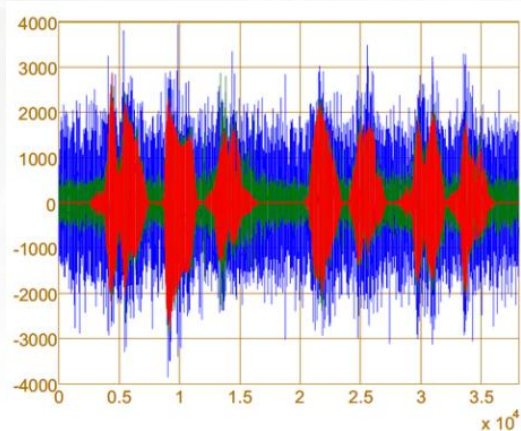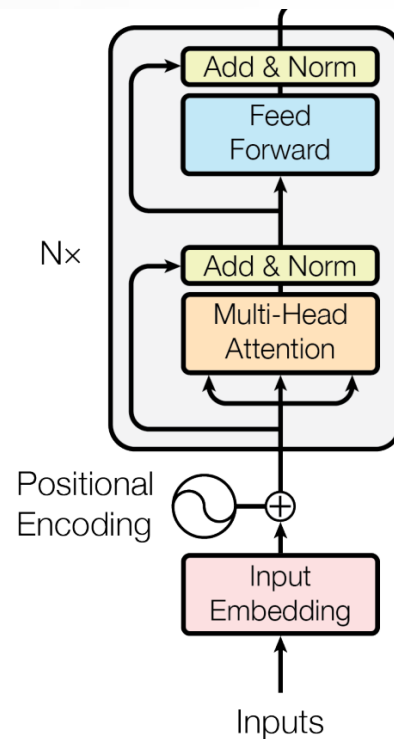| Feature Extraction Module | | Acoustic Neural Network | | Language Model | | Post Processing |
|:---:|:---:|:---:|:---:|:---:|:---:|:---:|
| Denoise&VAD、MFCC ... | → | Conformer、TDNN-LSTM... | → | N-gram、Neural Network... | → | Text Filtering Text Classification... |

# Challenge



THE SINITIC LANGUAGES

| | | |
|---|---|---|
| Mandarin | 836 million (worldwide) |
| Jin | 45 million(usu. grouped with Mandarin) |
| Wu | 77 million |
| Hui | 3.2 million(usu. grouped with Wu) |
| Gan | 31 million |
| Xiang | 36 million |
| Min | 60 million(incl. Taiwanese) |
| Hakka | 34 million (worldwide) |
| Yue | 71 million (worldwide) |
| Ping | 2 million(usu. grouped with Yue) |

Additive noise

Multiplicative noise

Reflective Sound Paths

Direct Sound Path

■ Negtive ■ Positive

**Diverse Accents**

**Multiple Devices & OS**

**Noisy Environment**

**Negative Samples Sparsity**

# Solution

| Massive Data Collecting | → | Valid Data Mining | → | Robust Acoustic Model Training |
|---|---|---|---|---|

**100,000+** hours
**accents/devices** coverage…

Closed-loop
Incremental…

# Valid Data Mining



Manually Labeled Data → (Training) → Initial Model ← (Update) ← New Model

Manually Labeled Data / Massive Audio Data → (Evaluation) → Initial Model

Test Data → Initial Model, New Model

Massive Audio Data → (Selection) → High Confidence Data → (Manual QA) → Incremental Data → (Training) → New Model

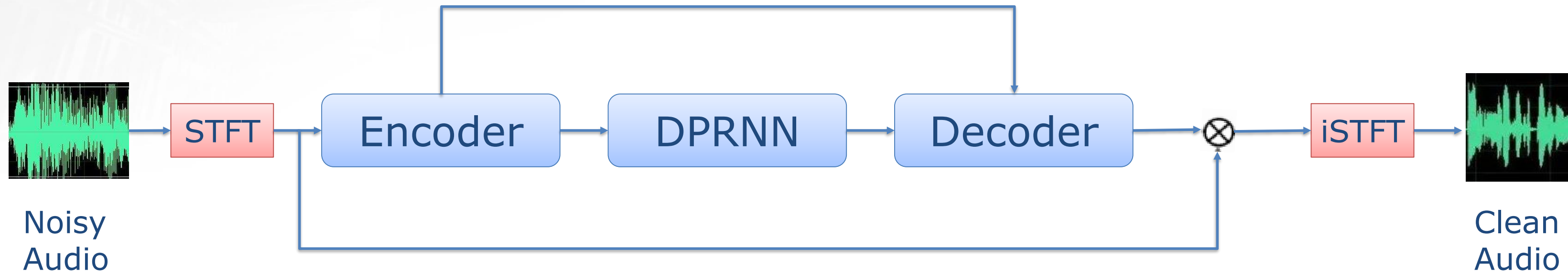# Noise Reduction

- Simulation data in 3000+ scenes
- Following DPCRN Network[1] but more lightweight



[1] Le X, Chen H, Chen K, et al. DPCRN: Dual-Path Convolution Recurrent Network for Single Channel Speech Enhancement.

# Keyword Enhancement



b*stard     <eps>

f*ck     off

sick     <eps>

Offensive words subgraph

You

make     me

What     happy

...

b*tch     <eps>

stupid     jerk

...

Main decoding graph

New offensive words subgraph

# Outline

- Introduction
- Image Sub-system
- Audio Sub-system
- **Text Sub-system**
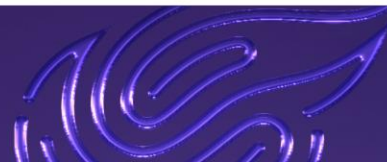- Application

# Text Filtering

www.163.c0m
USD|RMB|EUR
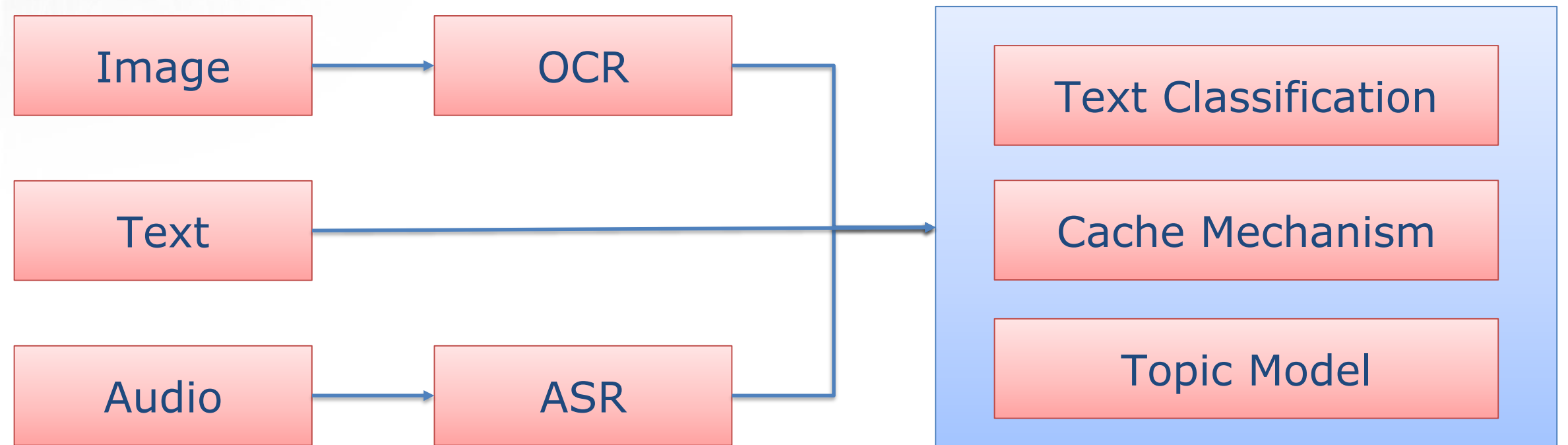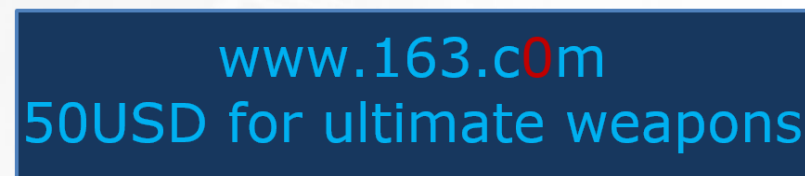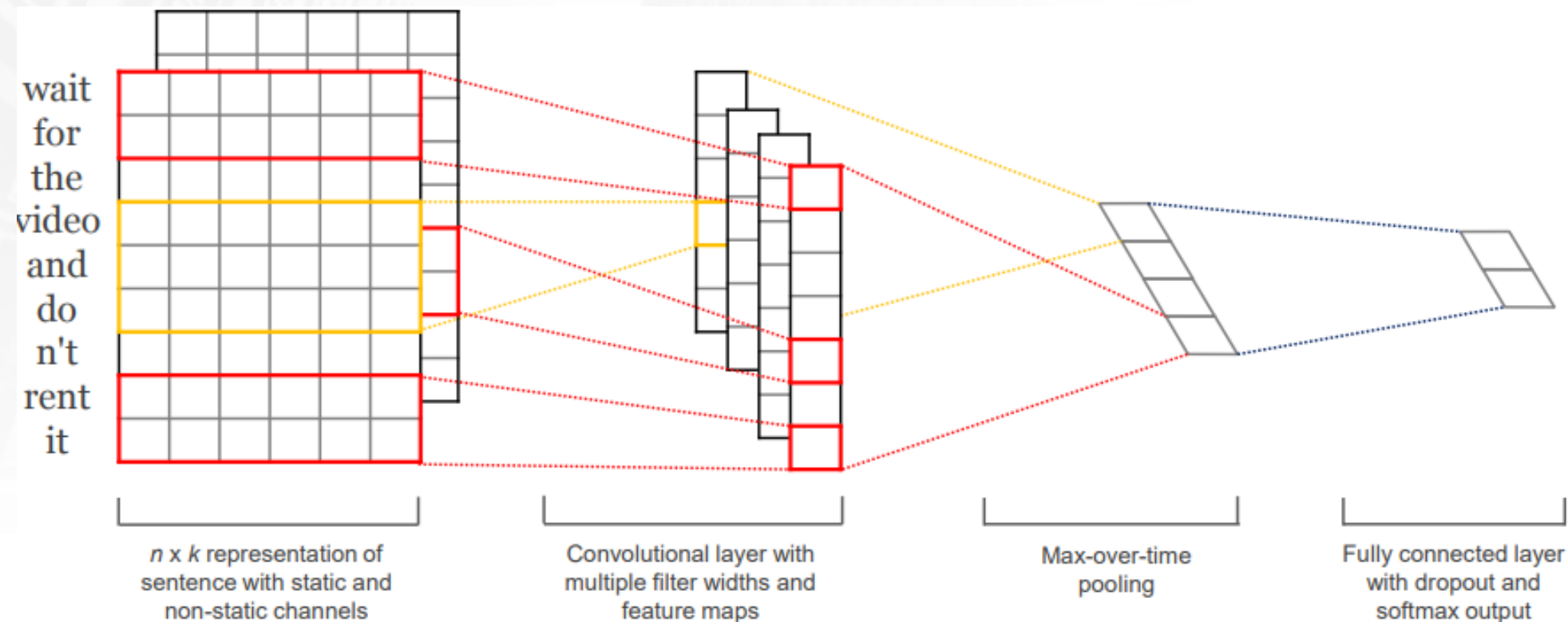
Regular
Expression

MACHINE
LEARNING

Machine
Learning
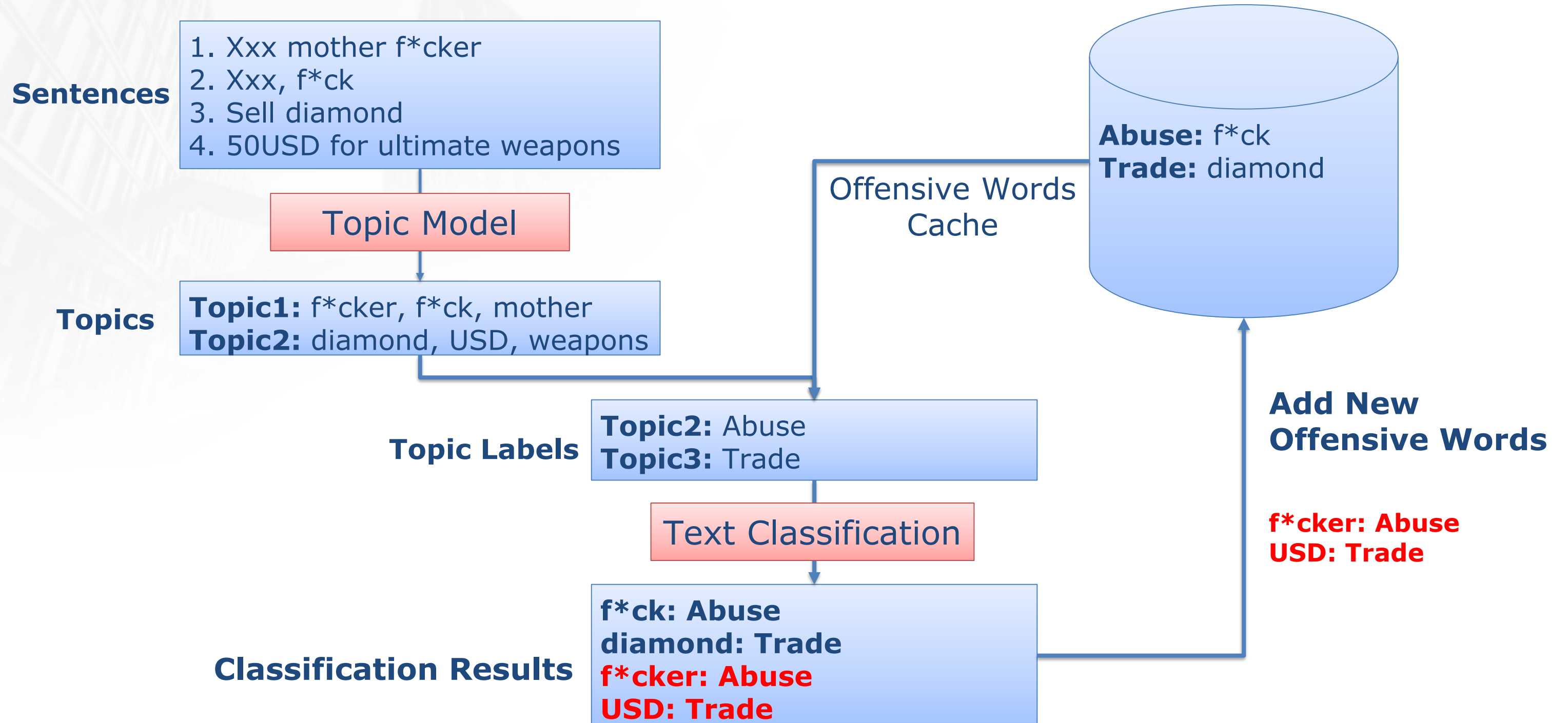
# Architecture

# Text Classification



Text CNN

- Use **multiple filters** with varying window sizes to obtain **multiple features**
- Quite simple
- High inference speed

[1] Chen Y. Convolutional neural network for sentence classification.

# Cache Mechanism

**Sentences**
1. Xxx mother f*cker
2. Xxx, f*ck
3. Sell diamond
4. 50USD for ultimate weapons

Topic Model

**Topics**
**Topic1:** f*cker, f*ck, mother
**Topic2:** diamond, USD, weapons

Offensive Words Cache

**Abuse:** f*ck
**Trade:** diamond

**Topic Labels**
**Topic2:** Abuse
**Topic3:** Trade

Text Classification

**Classification Results**
**f*ck: Abuse**
**diamond: Trade**
**f*cker: Abuse**
**USD: Trade**

**Add New Offensive Words**

**f*cker: Abuse**
**USD: Trade**

# Outline

- Introduction
- Image Sub-system
- Audio Sub-system
- Text Sub-system
- **Application**

# Architecture

HTTPS
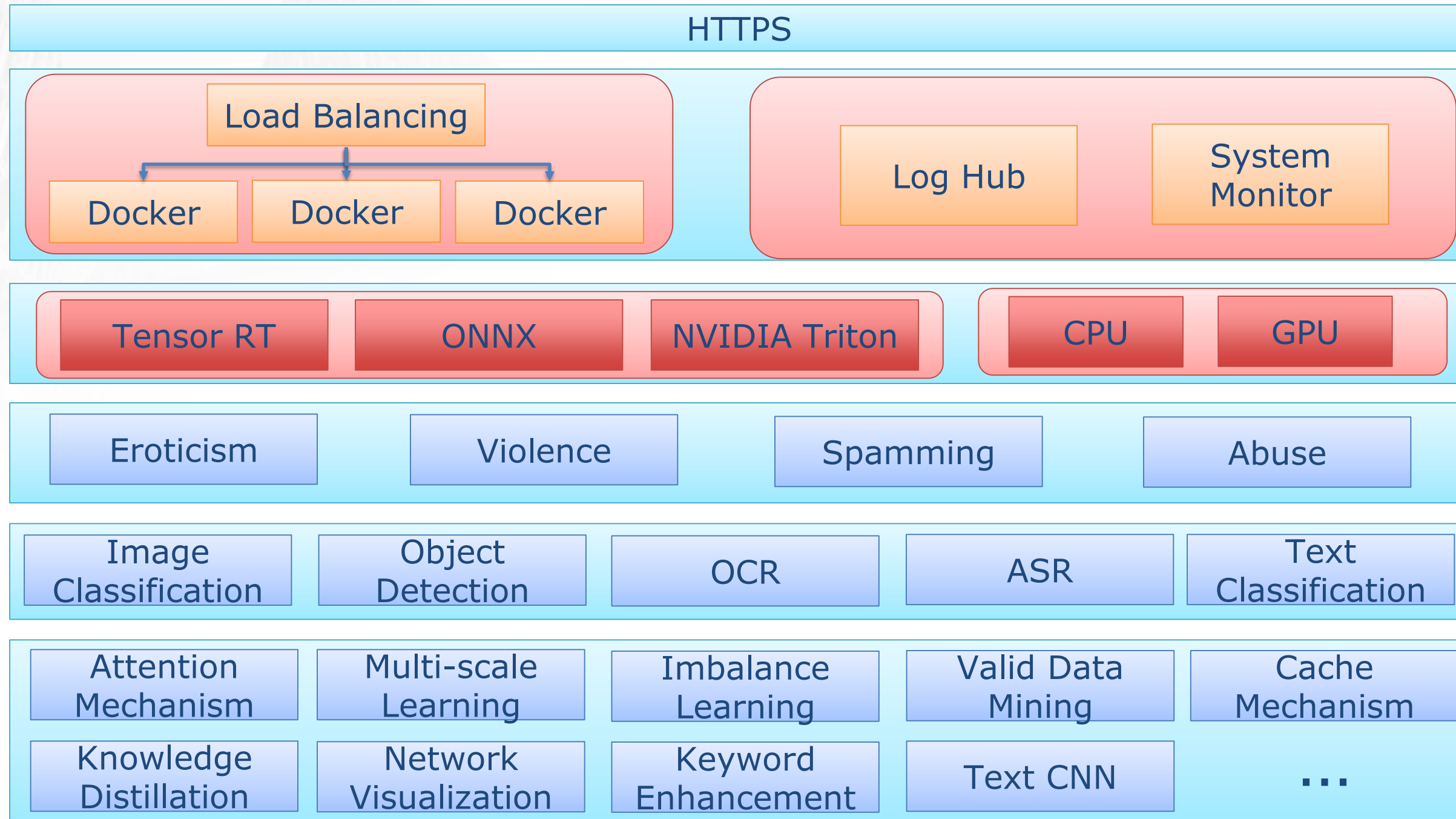
Load Balancing

Docker    Docker    Docker

Log Hub    System Monitor

| Tensor RT | ONNX | NVIDIA Triton | | CPU | GPU |

| Eroticism | Violence | Spamming | Abuse |

| Image Classification | Object Detection | OCR | ASR | Text Classification |

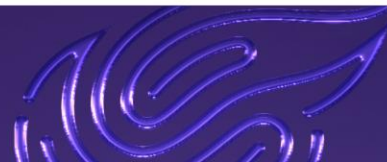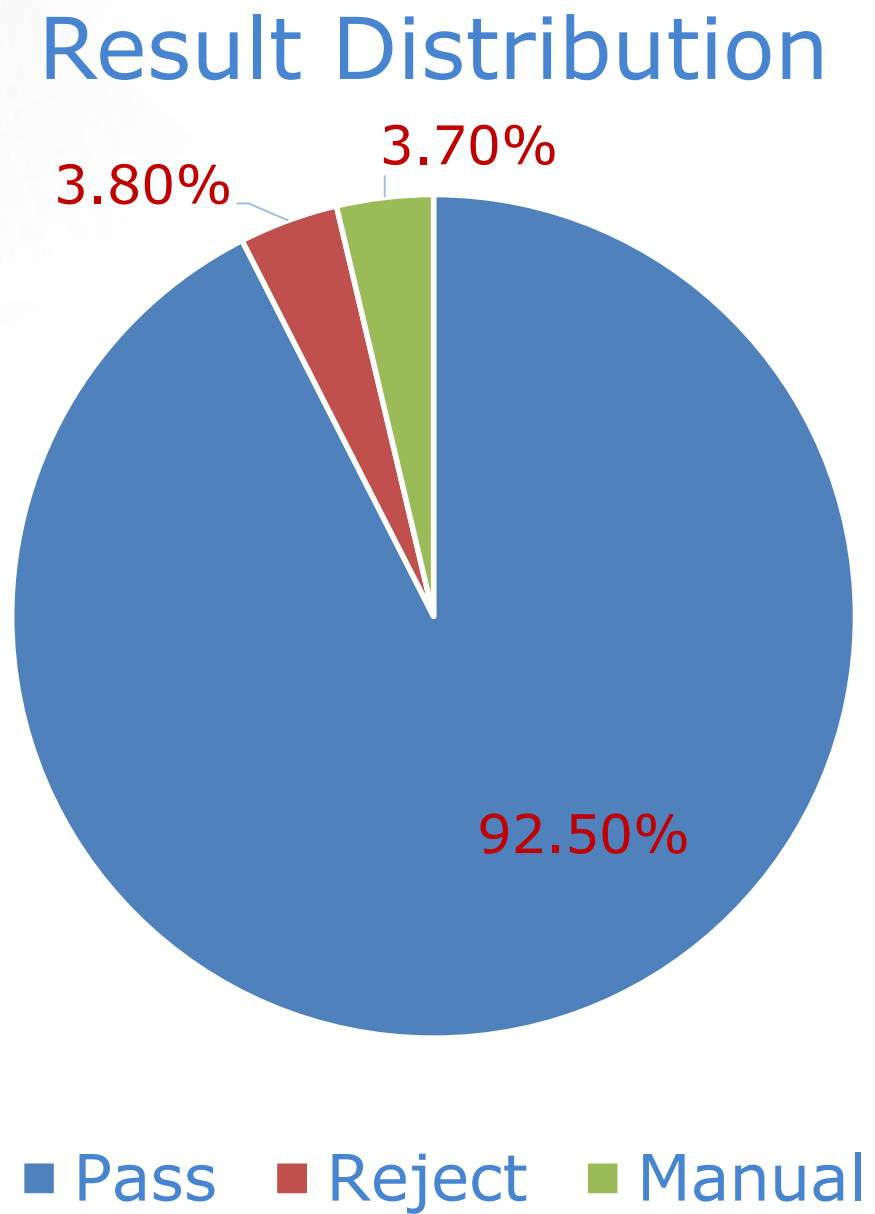| Attention Mechanism | Multi-scale Learning | Imbalance Learning | Valid Data Mining | Cache Mechanism |
| Knowledge Distillation | Network Visualization | Keyword Enhancement | Text CNN | ... |

# Application

- **High precision:** **98%** for image and **90%** for audio
- **High performance:** In total, more than **3,200,000,000** images and **73,000,000** hours of audio data are processed
- Has been running **stably for several years**
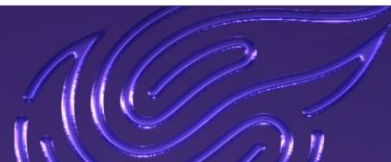- Has used in **almost all products** of NetEase Games

# Result Distribution

# Takeaways

- When build a multi-modal moderation system
  - **Data is critical.** We have introduced the methods for data collection, cleaning, mining and improvement.
  - **Model update and optimization never stop**. Some of the methods we adopted are shown and can be used as a reference.