# A **Comparison** of
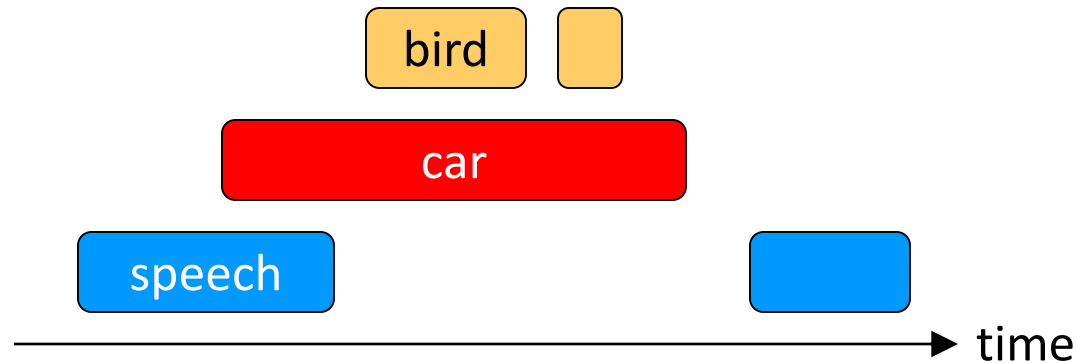# **Five Multiple Instance Learning** <span style="color:red">Pooling Functions</span>
# for <span style="color:red">Sound Event Detection</span> with **Weak Labeling**

Yun Wang, **Juncheng L**i, Florian Metze

May 14, 2019

# Sound Event Detection

- Detection = audio tagging + <span style="color:red">localization</span>



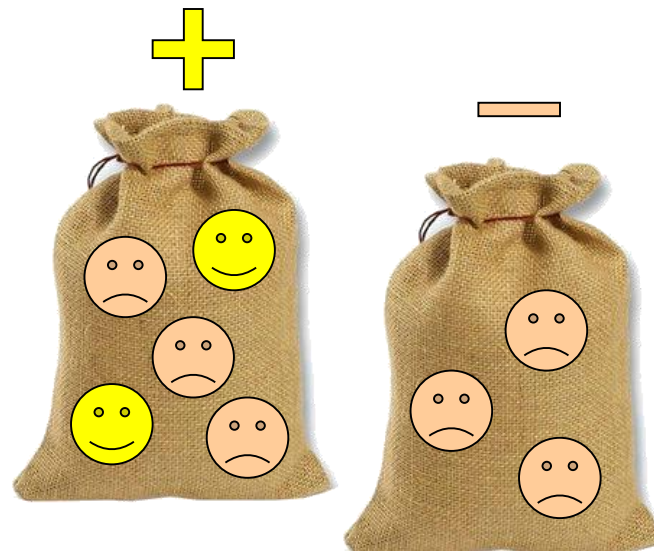- Strong labeling is expensive to obtain

# Sound Event Detection
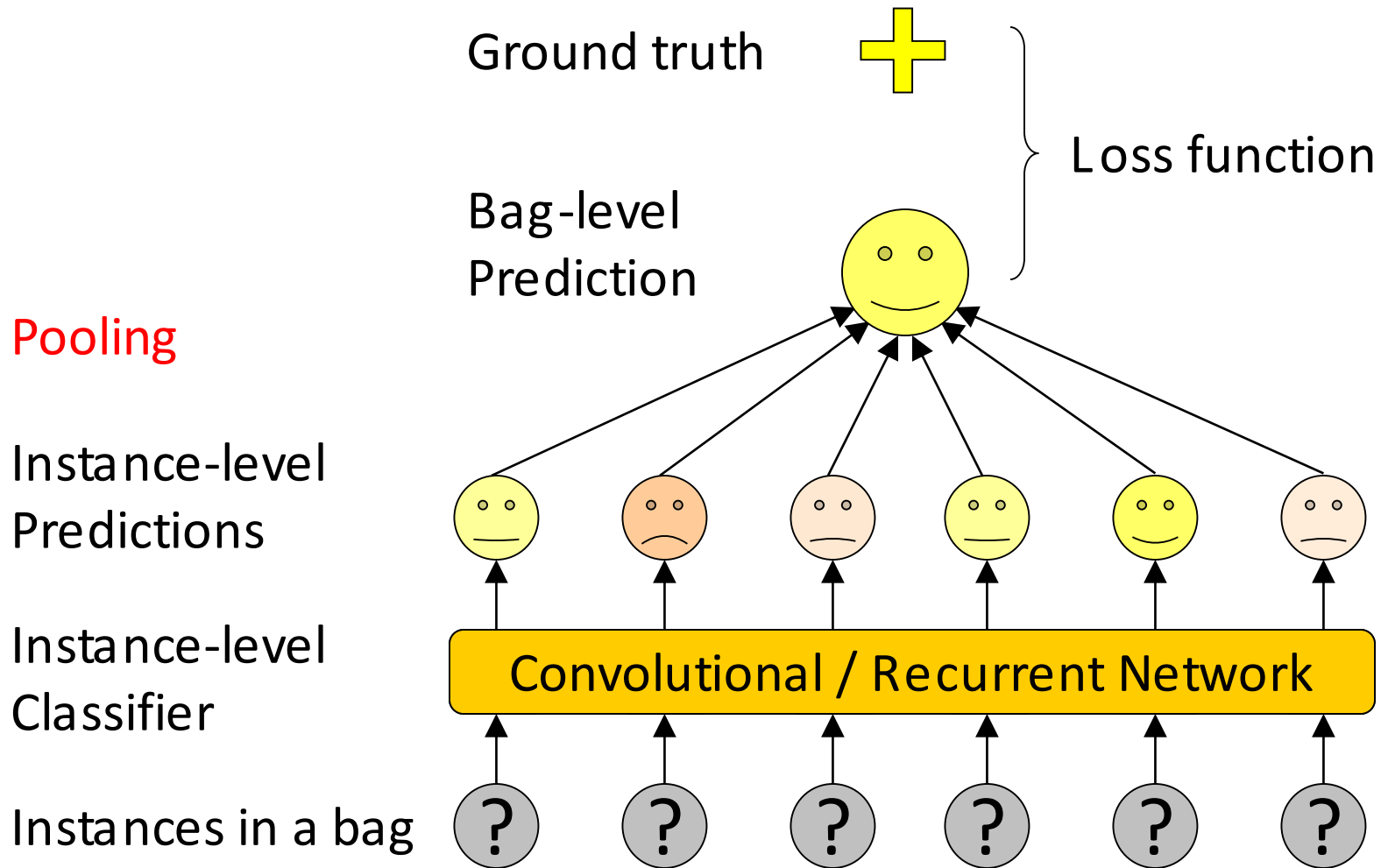
- Train with <span style="color:red">weak labeling</span>

- But still, we want both <span style="color:red">tagging</span> and <span style="color:red">localization</span> output

# Multiple Instance Learning

- SED with weak labeling is a <span style="color:red">Multiple Instance Learning</span> (MIL) problem
  - Bag is positive ⟺ any instance is positive
  - Recording = bag, frames = instances

# Multiple Instance Learning

Ground truth

Loss function

Bag-level Prediction

Pooling

Instance-level Predictions

Instance-level Classifier

Convolutional / Recurrent Network

Instances in a bag

5

# Pooling Functions

**Max pooling**

$$y = \max_i y_i$$

**Linear softmax**

$$y = \frac{\sum_i y_i^2}{\sum_i y_i}$$

**Exp. softmax**

$$y = \frac{\sum_i y_i \exp(y_i)}{\sum_i \exp(y_i)}$$

**Average pooling**

$$y = \frac{1}{n}\sum_i y_i$$

Weighted Average

One frame gets
all the weight

Larger probs
get larger weight

All frames get
equal weight

$$y = \frac{\sum_i y_i w_i}{\sum_i w_i}$$

Attention:
Learn the weights!

# Pooling Functions

- We found linear softmax best for localization!

$$y = \frac{\sum_i y_i^2}{\sum_i y_i} \qquad \frac{\partial y}{\partial y_i} = \frac{2y_i - y}{\sum_j y_j}$$

Positive when $y_i > y/2$

- When bag is positive:
  - $y_i$ gets away from $y/2$
  - Only boosts frames with $y_i > y/2$ – nice localization!
- When bag is negative:
  - $y_i$ approaches $y/2$ – finally converges to zero

# Pooling Functions

- **What's wrong with <span style="color:red">attention</span>?**

$$y = \frac{\sum_i y_i w_i}{\sum_i w_i} \qquad \frac{\partial y}{\partial y_i} = \frac{w_i}{\sum_j w_j} \qquad \frac{\partial y}{\partial w_i} = \frac{y_i - y}{\sum_j w_j}$$

<span style="background-color:red; color:yellow">Always positive</span>  <span style="background-color:red; color:yellow">Positive when $y_i > y$</span>

- **When bag is positive:**
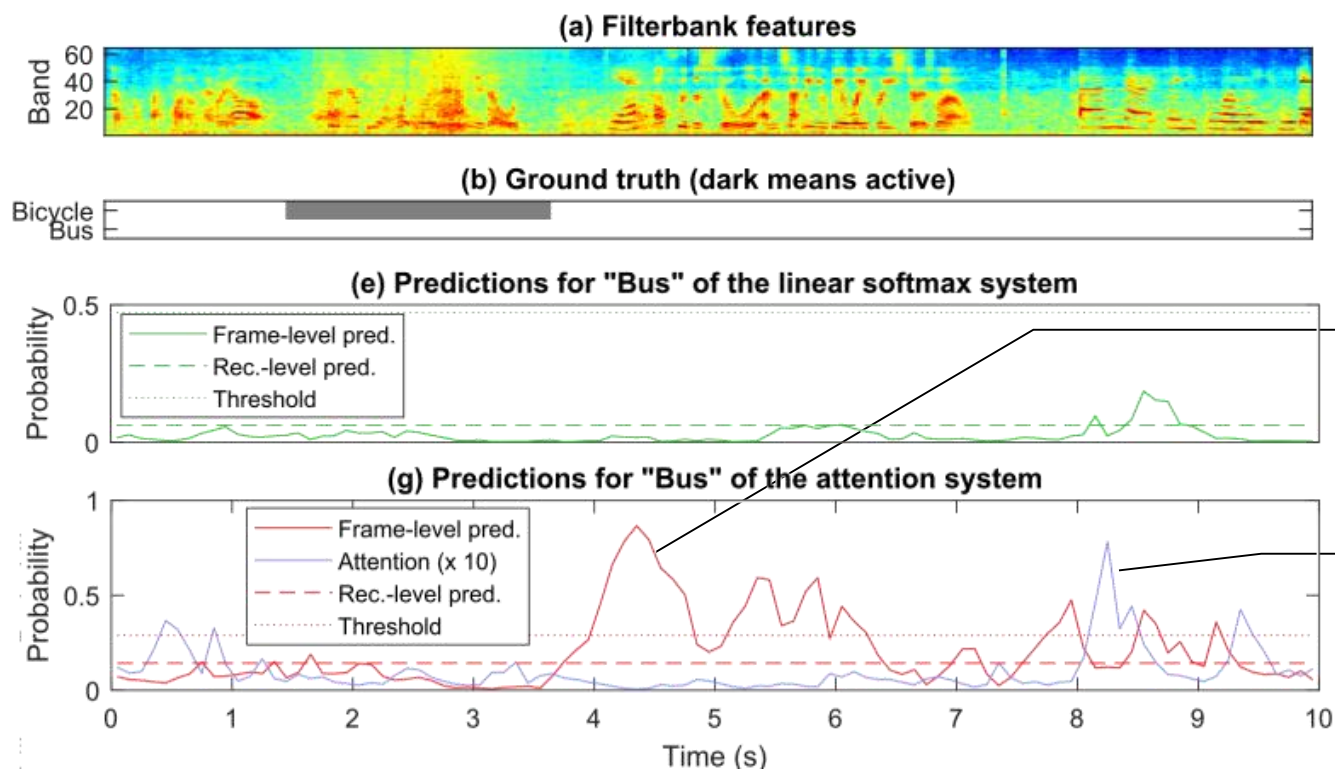  - ❑ All $y_i$ increase 🙂 attention focuses where $y_i > y$   🙂
- **When bag is negative:**
  - ❑ All $y_i$ decrease 🙂 attention focuses where $y_i < y$   😱
  - ❑ <span style="color:red">Smaller probs get larger weight!</span>

8

# Failure Mode of Attention



(a) Filterbank features

(b) Ground truth (dark means active)

(e) Predictions for "Bus" of the linear softmax system

(g) Predictions for "Bus" of the attention system

False positives in unattended regions

Attention focuses here

- Too many frame-level false positives

- Inconsistent recording-level and frame-level predictions
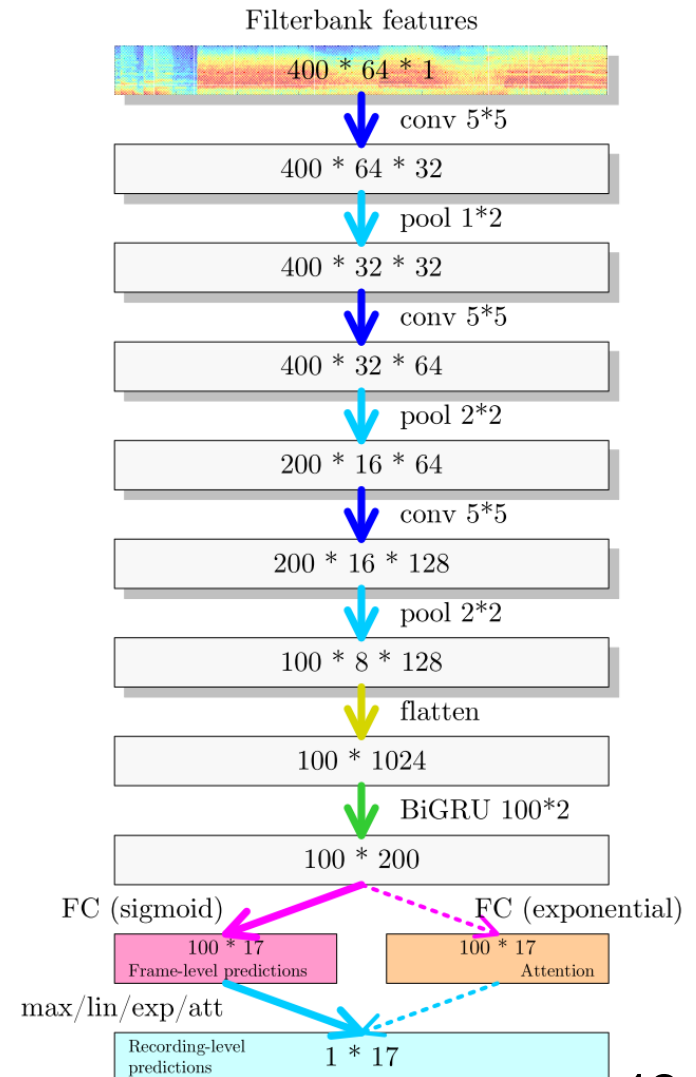
# EVALUATION I:
## DCASE 2017 Challenge, Task 4

# DCASE 2017: Task

- **17 event types**
  - Vehicles, warnings
- **Training data:**
  - ~50k recordings * 10 seconds each = ~140 hours
  - Weakly labeled
- **Test data:**
  - 488 recordings * 10 seconds each = ~1.4 h
  - Strongly labeled
- **Evaluation metrics:**
  - Tagging: F1
  - Localization: error rate & F1 on 1s segments

# DCASE 2017: Model

- Input:
  - Logmel features @ 40 Hz
- Structure:
  - 3 conv layers + 1 GRU layer
- Output:
  - Frame-level event probs at 10 Hz
  - For tagging: pooled globally into recording-level event probs
  - For localization: pooled over 1s segments

Filterbank features

400 * 64 * 1

conv 5*5

400 * 64 * 32

pool 1*2

400 * 32 * 32

conv 5*5

400 * 32 * 64

pool 2*2

200 * 16 * 64

conv 5*5

200 * 16 * 128

pool 2*2

100 * 8 * 128

flatten

100 * 1024

BiGRU 100*2

100 * 200

FC (sigmoid)        FC (exponential)

100 * 17
Frame-level predictions

100 * 17
Attention

max/lin/exp/att

Recording-level predictions        1 * 17

# DCASE 2017: Results

| Pooling Func | Tag F1 | Loc ER | Loc F1 | Loc #FN | Loc #FP |
|---|---|---|---|---|---|
| Max | 45.3 | 84.7 | 35.4 | 3,154 | 1,253 |
| Linear softmax | 49.5 | 84.3 | 43.7 | 2,528 | 2,187 |
| Attention | 49.2 | 102.5 | 40.1 | 2,434 | 3,309 |

- Max: too many false negatives (FNs) hurt F1
- Attention: too many false positives (FPs) hurt ER
- Linear softmax: balanced FNs and FPs

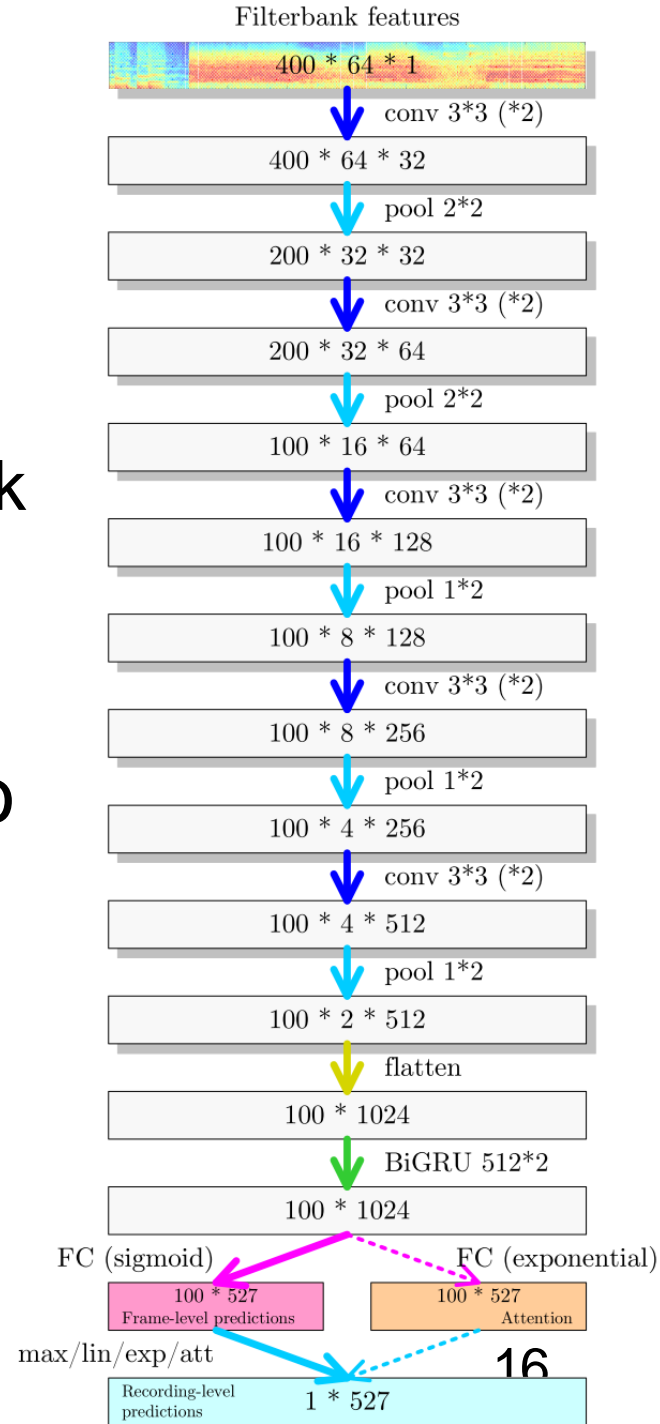# EVALUATION II:
## Google Audio Set

# Audio Set: Task

- **Data:**
  - 527 event types (include the 17 events of DCASE)
  - Weakly labeled
  - Training: ~2M recordings * 10s = 8 months
  - Test: ~20k recordings * 10s = 56 hours
- **Evaluation metrics:**
  - Audio Set only measures <span style="color:red">tagging</span>
    - MAP, MAUC, d'
  - Reuse DCASE data & metrics for <span style="color:red">tagging</span> & <span style="color:red">localization</span>
    - Tag F1, Loc ER, Loc F1 over 1s segments

# Audio Set: Model

- **TALNet**:
  - Tagging and Localization Network
  - 10 conv layers, 1 GRU layer
  - Same input & output as before
- No fine-tuning when applied to DCASE data

Filterbank features

400 * 64 * 1

conv 3*3 (*2)

400 * 64 * 32

pool 2*2

200 * 32 * 32

conv 3*3 (*2)

200 * 32 * 64

pool 2*2

100 * 16 * 64

conv 3*3 (*2)

100 * 16 * 128

pool 1*2

100 * 8 * 128

conv 3*3 (*2)

100 * 8 * 256

pool 1*2

100 * 4 * 256

conv 3*3 (*2)

100 * 4 * 512

pool 1*2

100 * 2 * 512

flatten

100 * 1024

BiGRU 512*2

100 * 1024

FC (sigmoid)                    FC (exponential)

100 * 527
Frame-level predictions          100 * 527
                                 Attention

max/lin/exp/att                                    16

Recording-level predictions          1 * 527

# Audio Set: Result 1/3

| Group | System | No. of Training Recs. | Audio Set | | | DCASE 2017 | | |
| | | | MAP | MAUC | d' | Task A | Task B | |
| | | | | | | F1 | ER | F1 |
| TALNet (Sec. 3.3) | Max pooling | 2M | 0.351 | 0.961 | 2.497 | 52.6 | 81.5 | 42.2 |
| | Average pooling | | 0.361 | **0.966** | 2.574 | **53.8** | 101.8 | **46.8** |
| | Linear softmax | | 0.359 | **0.966** | **2.575** | 52.3 | **78.9** | 45.4 |
| | Exp. softmax | | **0.362** | 0.965 | 2.554 | 52.3 | 89.2 | 46.2 |
| | Attention | | 0.354 | 0.963 | 2.531 | 51.4 | 92.0 | 45.5 |

- **TALNet works out of the box on DCASE**

- **Linear softmax is best for localization**

  - And good enough for tagging

17

# Audio Set: Result 2/3

| Group | System | No. of Training Recs. | Audio Set | | | DCASE 2017 | | |
|---|---|---|---|---|---|---|---|---|
| | | | MAP | MAUC | d' | Task A | Task B | |
| | | | | | | F1 | ER | F1 |
| TALNet (Sec. 3.3) | Max pooling | 2M | 0.351 | 0.961 | 2.497 | 52.6 | 81.5 | 42.2 |
| | Average pooling | | 0.361 | **0.966** | 2.574 | **53.8** | 101.8 | **46.8** |
| | Linear softmax | | 0.359 | **0.966** | **2.575** | 52.3 | **78.9** | 45.4 |
| | Exp. softmax | | **0.362** | 0.965 | 2.554 | 52.3 | 89.2 | 46.2 |
| | Attention | | 0.354 | 0.963 | 2.531 | 51.4 | 92.0 | 45.5 |
| Literature | Hershey [71, 15] | 1M | 0.314 | 0.959 | 2.452 | | | |
| | Kumar [128] | 22k | 0.213 | 0.927 | | | | |
| | Shah [48] | 22k | 0.229 | 0.927 | | | | |
| | Wu [131] | 22k | | 0.927 | | | | |
| | Kong [54] | 2M | 0.327 | 0.965 | 2.558 | | | |
| | Yu [55] | 2M | **0.360** | **0.970** | **2.660** | | | |
| | Chen [56] | 600k | 0.316 | | | | | |
| | Chou [57] | 1M | 0.327 | 0.951 | | | | |

- **TALNet closely matches state of the art on tagging**
  - Yu's system uses multi-level attention and can't do localization!
- Amount of training data matters!

# Audio Set: Result 3/3

| Group | System | No. of Training Recs. | Audio Set | | | DCASE 2017 | | |
|---|---|---|---|---|---|---|---|---|
| | | | MAP | MAUC | d' | Task A | Task B | |
| | | | | | | F1 | ER | F1 |
| TALNet (Sec. 3.3) | Max pooling | 2M | 0.351 | 0.961 | 2.497 | 52.6 | 81.5 | 42.2 |
| | Average pooling | | 0.361 | **0.966** | 2.574 | **53.8** | 101.8 | **46.8** |
| | Linear softmax | | 0.359 | **0.966** | **2.575** | 52.3 | **78.9** | 45.4 |
| | Exp. softmax | | **0.362** | 0.965 | 2.554 | 52.3 | 89.2 | 46.2 |
| | Attention | | 0.354 | 0.963 | 2.531 | 51.4 | 92.0 | 45.5 |
| DCASE only (Sec. 3.2.3) | Max pooling | 50k | | | | 45.3 | 84.7 | 35.4 |
| | Average pooling | | | | | **50.0** | 105.9 | 41.3 |
| | Linear softmax | | | | | 49.5 | **84.3** | **43.7** |
| | Exp. softmax | | | | | 48.5 | 100.6 | 42.8 |
| | Attention | | | | | 49.2 | 102.5 | 40.1 |

- **Adding more data helps the 17 DCASE events**
  - Even though most of it belongs to 510 other events

# Summary

- Linear softmax is the best for localization
  - Better than max: unobstructed gradient flow
  - Better than attention:
    - Balanced false negatives and false positives
    - Consistent frame-level & recording-level predictions
- We built TALNet
  - First simultaneous audio tagging and localization
  - Closely matches state of the art on Audio Set
  - Good performance on DCASE 2017 out of the box
- Future work
  - Attention pooling with monotonicity constraint?

# Thanks!

# Questions?