

Weakly-Supervised Temporal Localization via Occurrence Count Learning



Julien Schroeter

schroeterj1@cardiff.ac.uk

Dr Kirill Sidorov

Prof David Marshall

CONTEXT

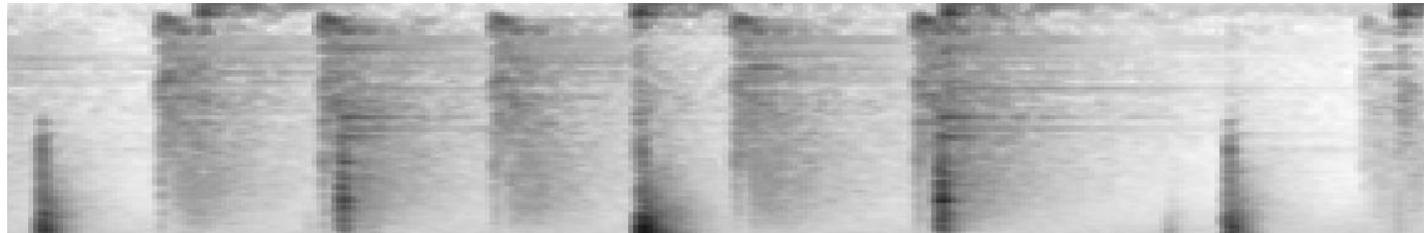
Temporal Localization



CONTEXT

Temporal Localization

Input Data



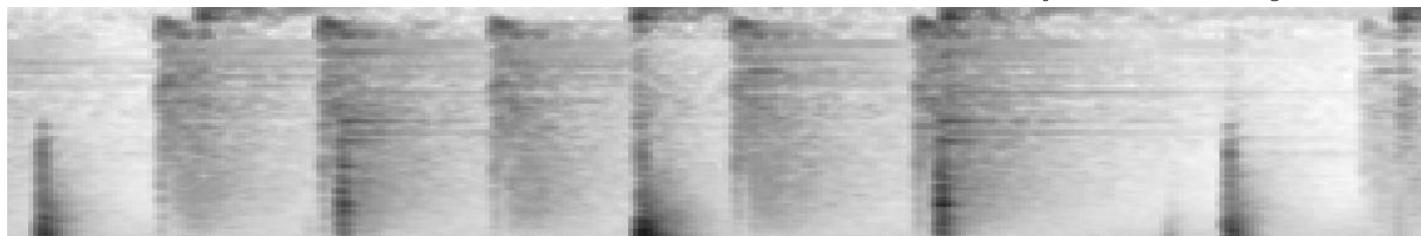
Temporal Sequence



CONTEXT

Temporal Localization

Input Data



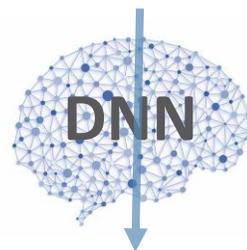
Temporal Sequence

Target



Impulse-like Events

Precise Localization





CONTEXT

Temporal Localization

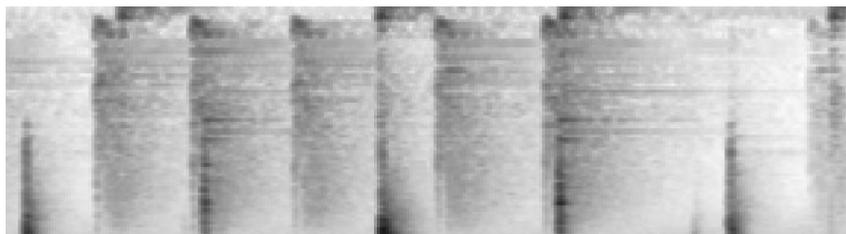


———— Fully-Supervised ————

CONTEXT

Temporal Localization

Input Data

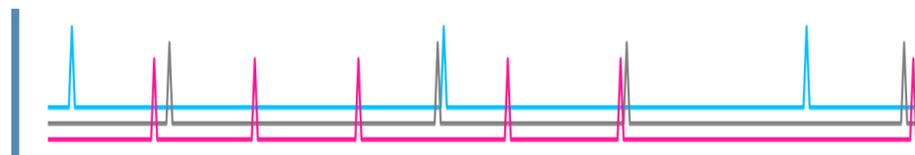


———— Fully-Supervised ————

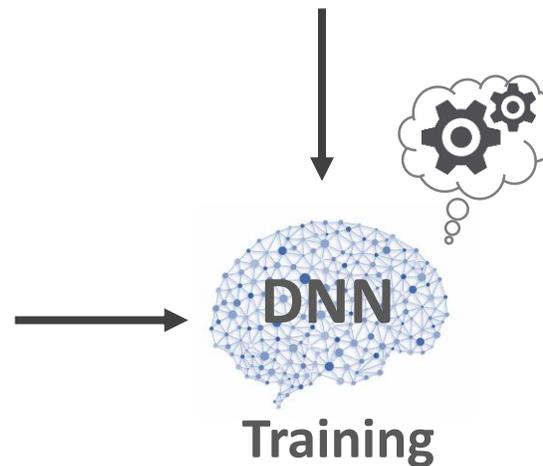
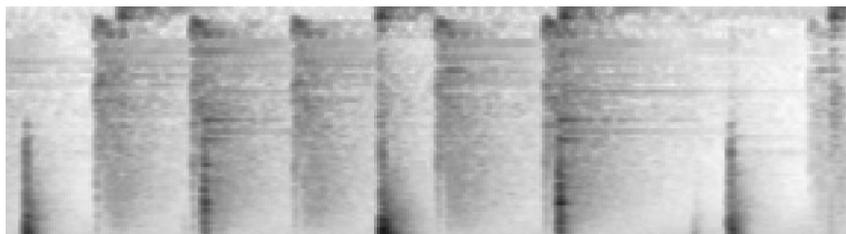
CONTEXT

Temporal Localization

Training Labels



Input Data



———— Fully-Supervised ————

OBJECTIVE

Weakening the annotation requirement

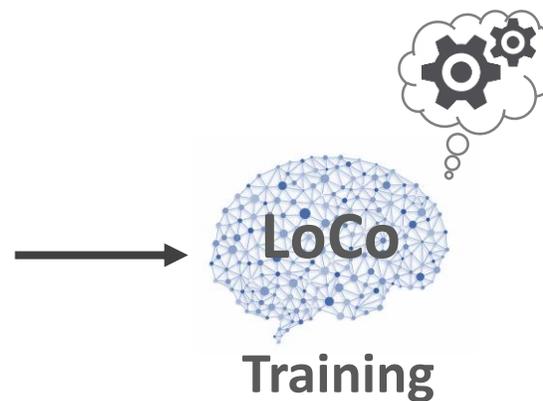
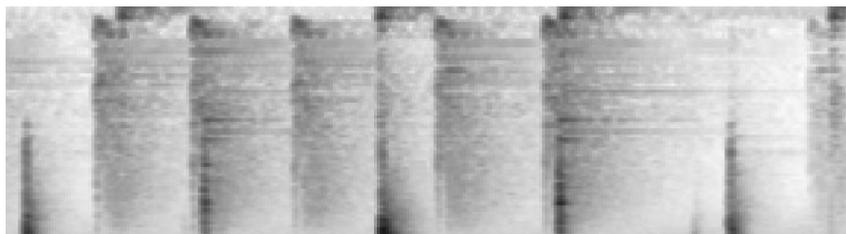


————— Our approach —————

OBJECTIVE

Weakening the annotation requirement

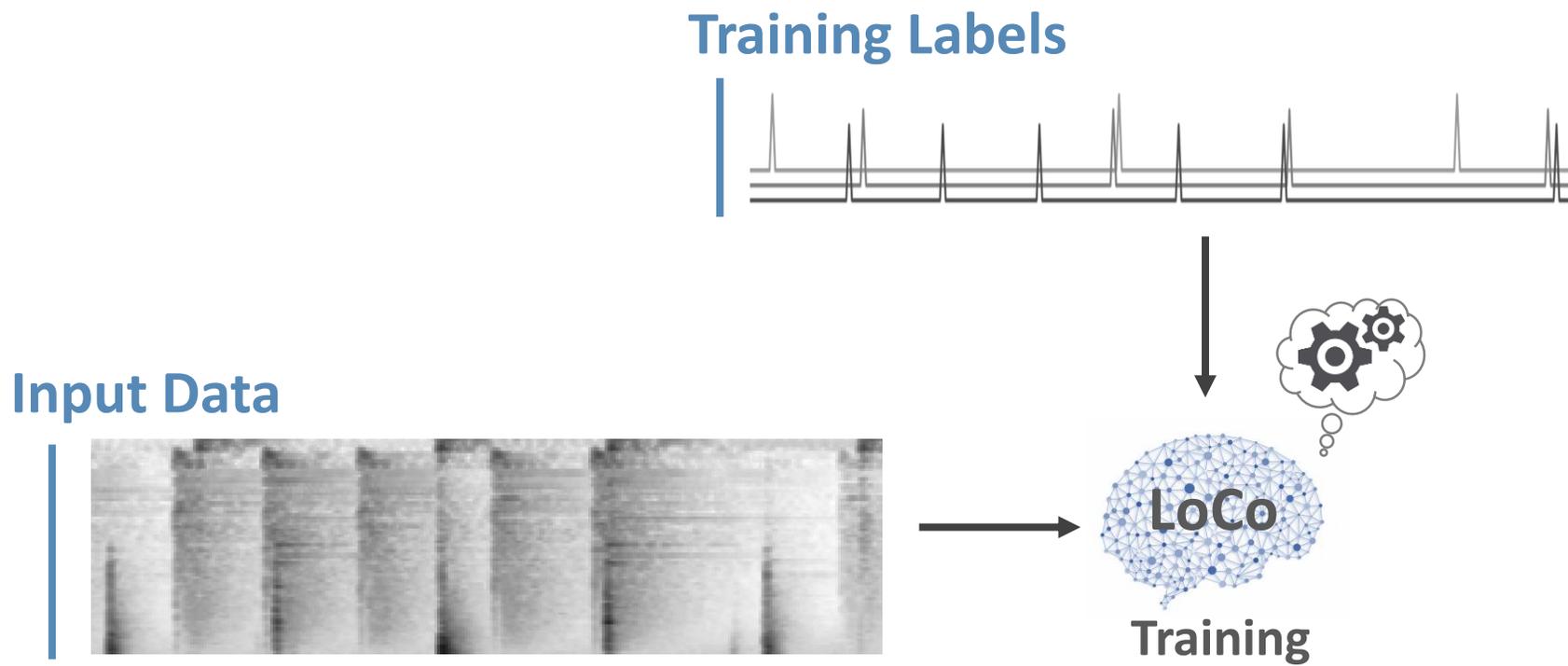
Input Data



Our approach

OBJECTIVE

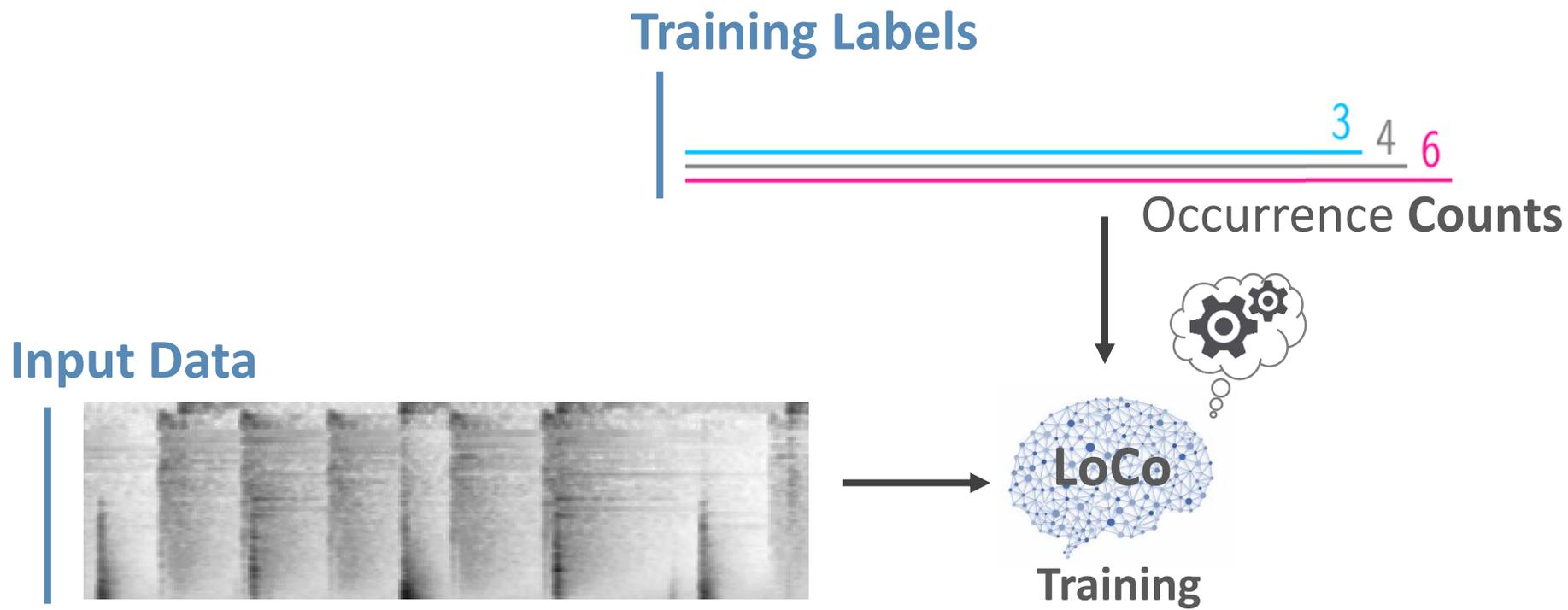
Weakening the annotation requirement



————— Our approach —————

OBJECTIVE

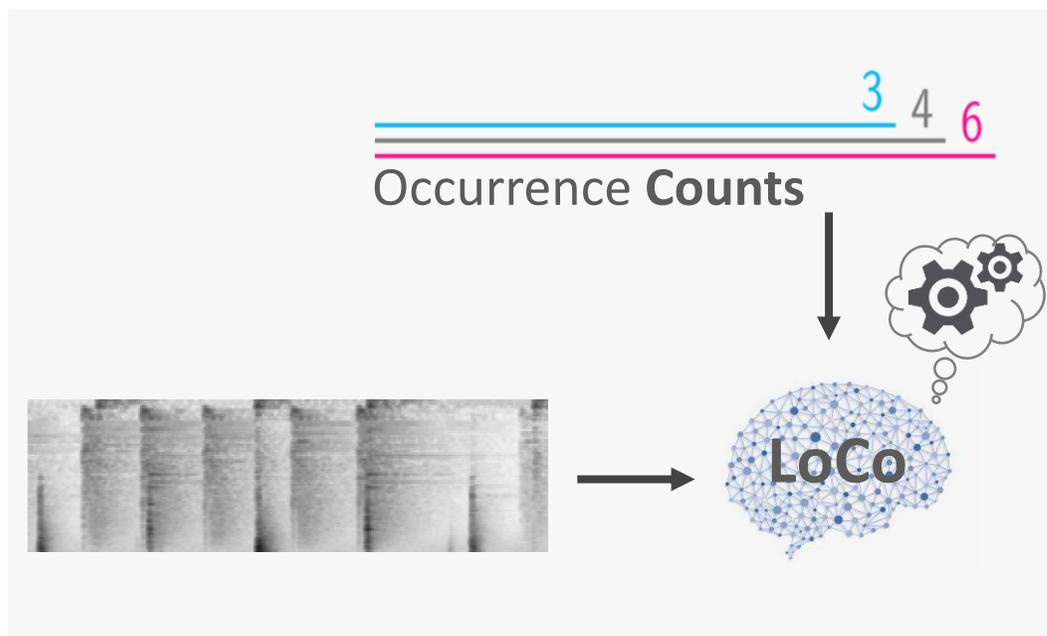
Weakening the annotation requirement



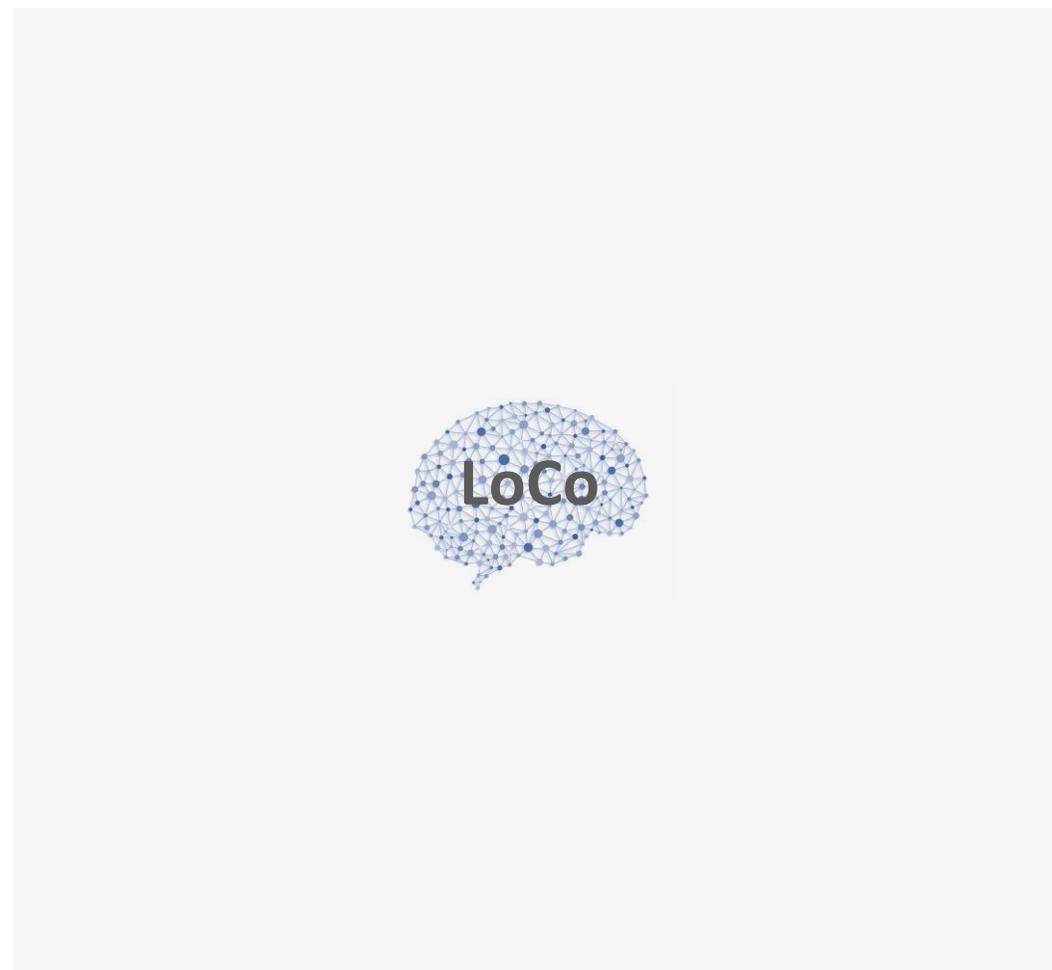
————— Our approach —————

OBJECTIVE

Weakening the annotation requirement



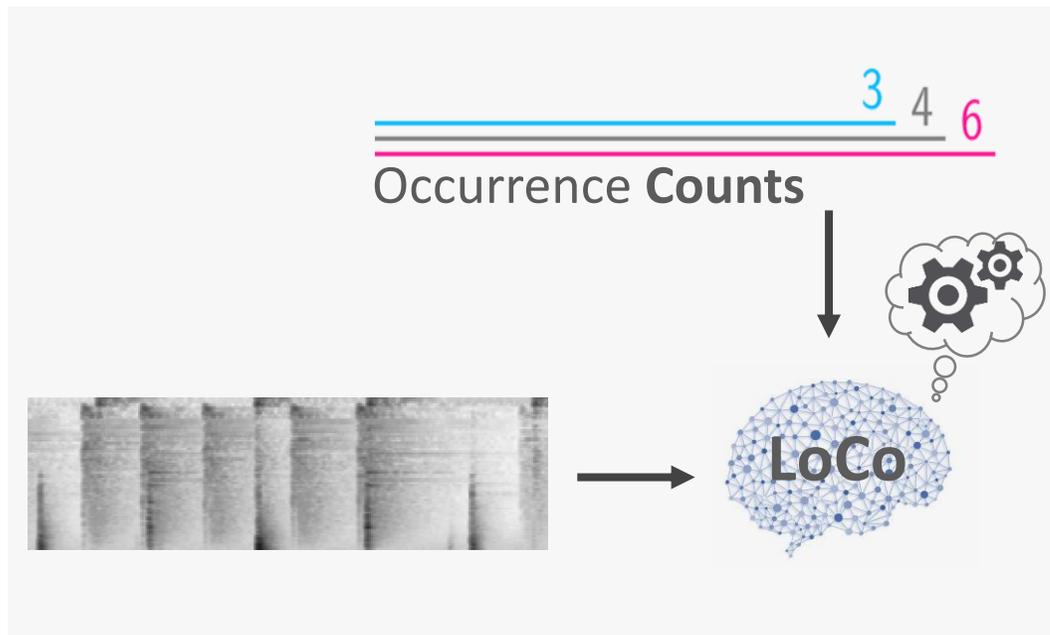
Training



Inference

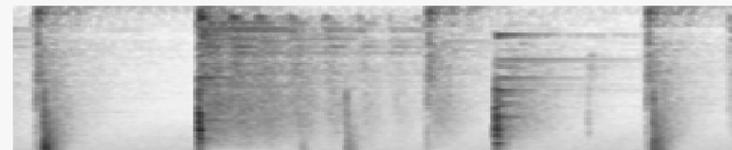
OBJECTIVE

Weakening the annotation requirement



Training

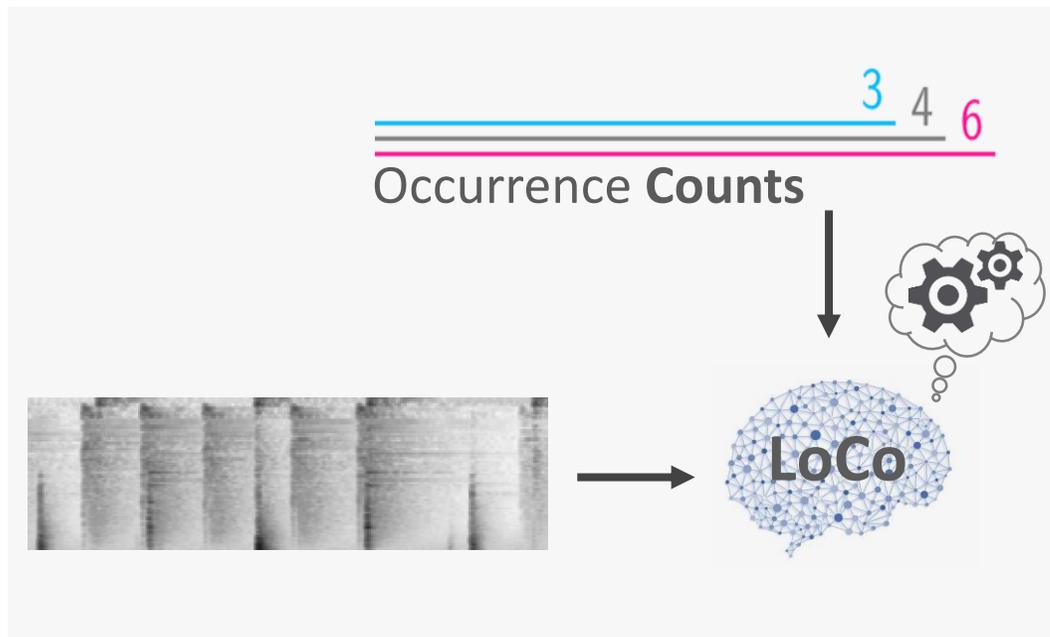
Input



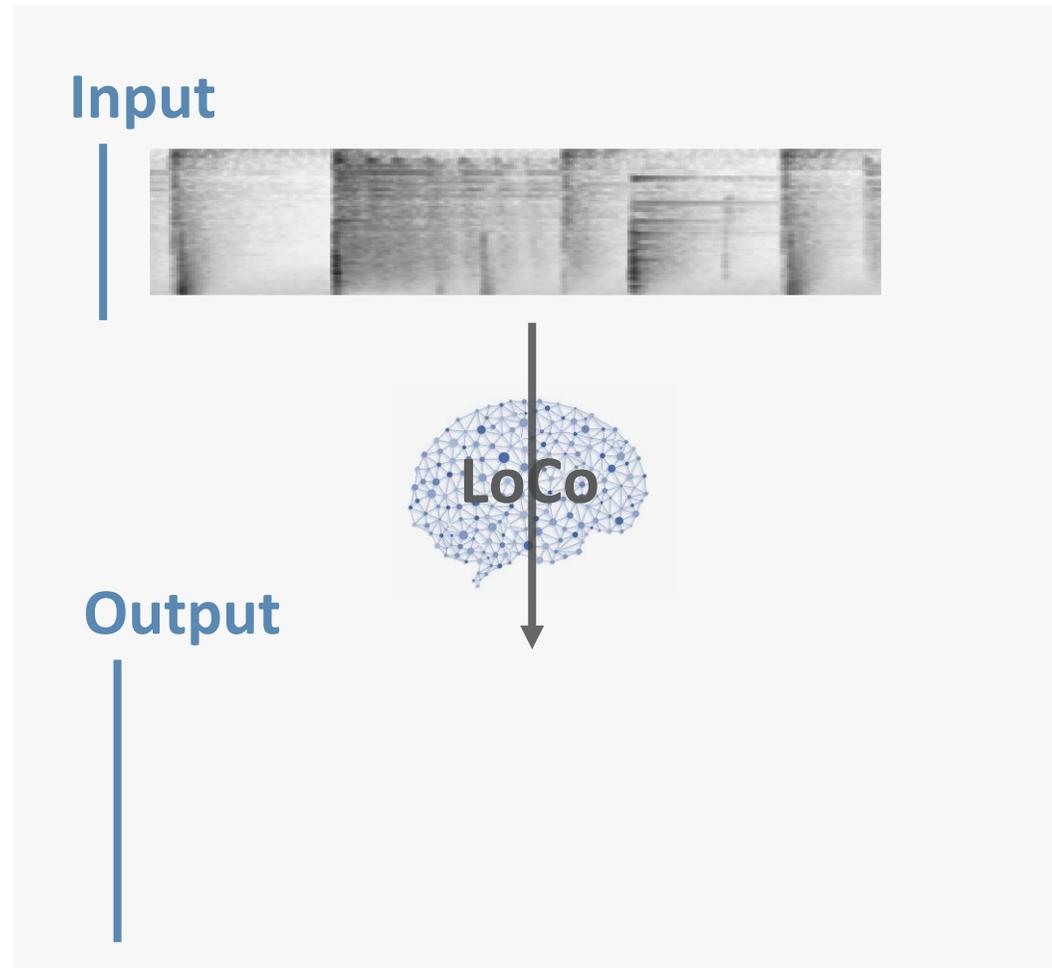
Inference

OBJECTIVE

Weakening the annotation requirement



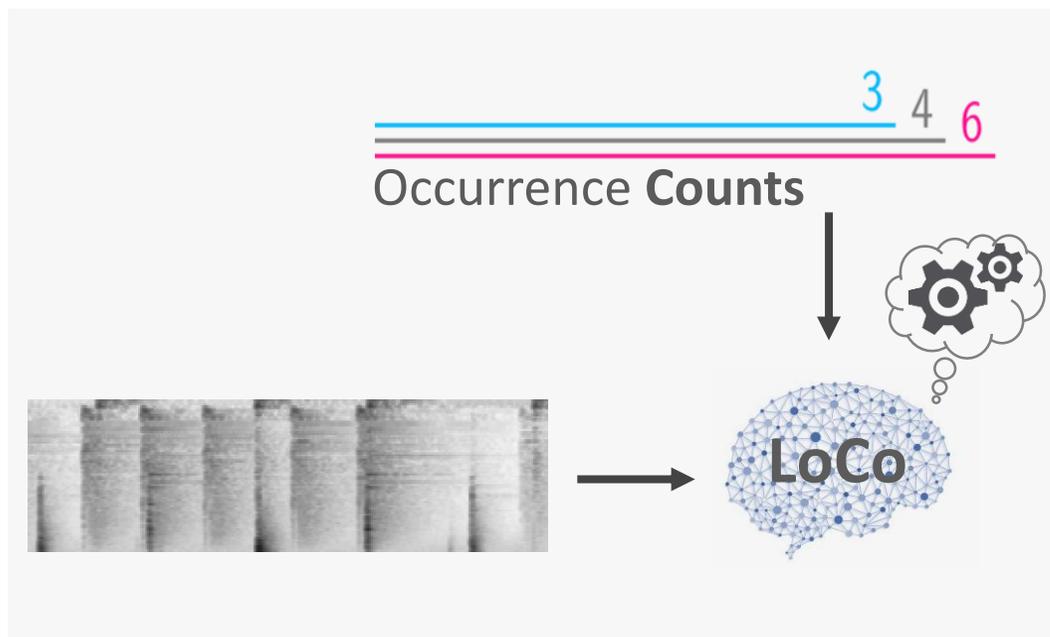
Training



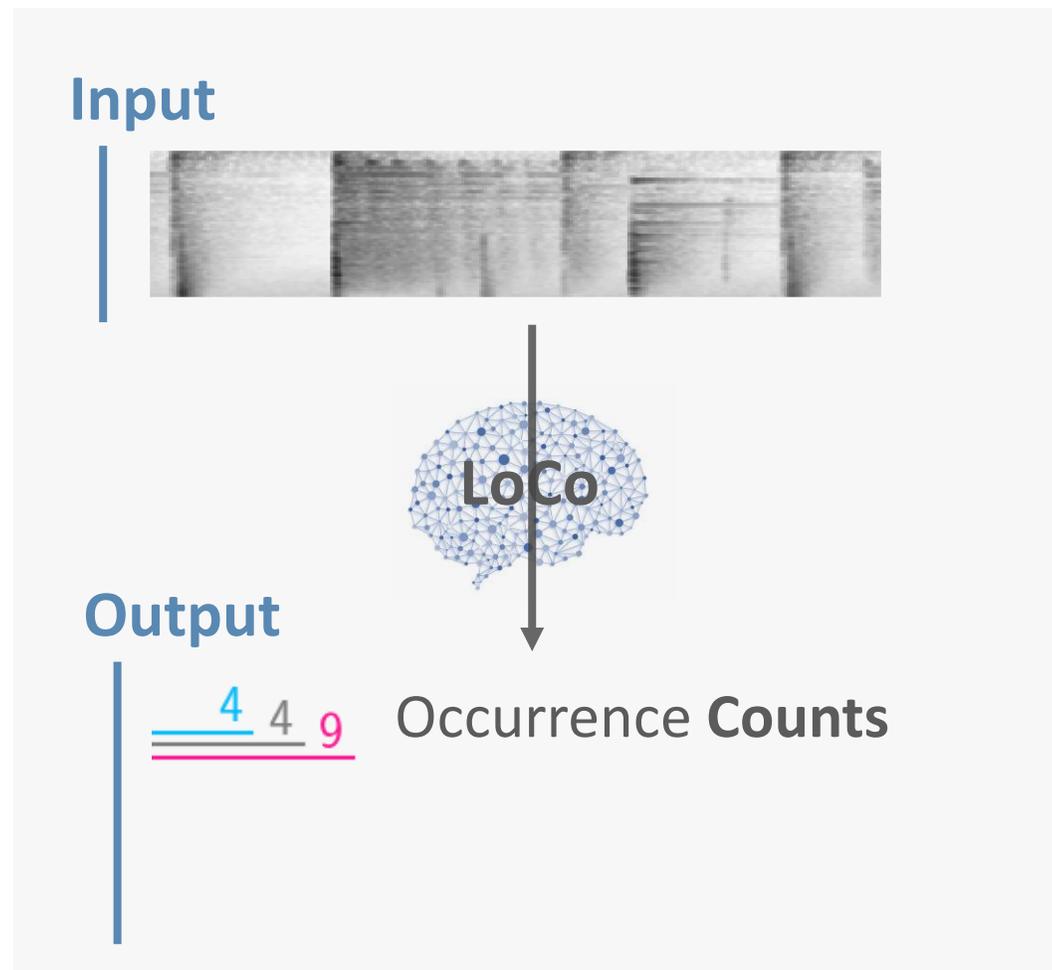
Inference

OBJECTIVE

Weakening the annotation requirement



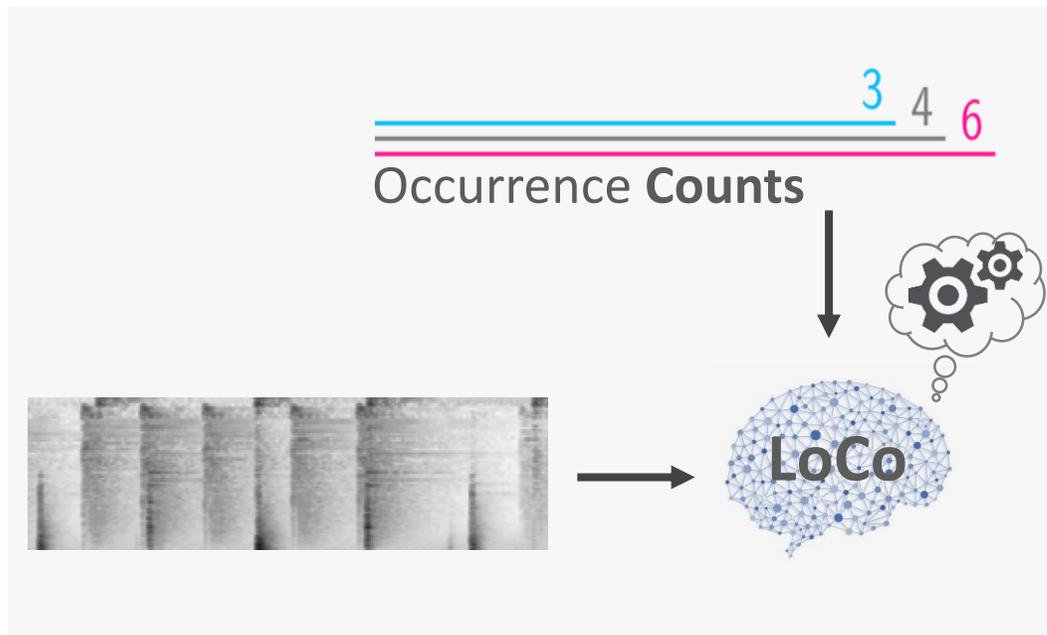
Training



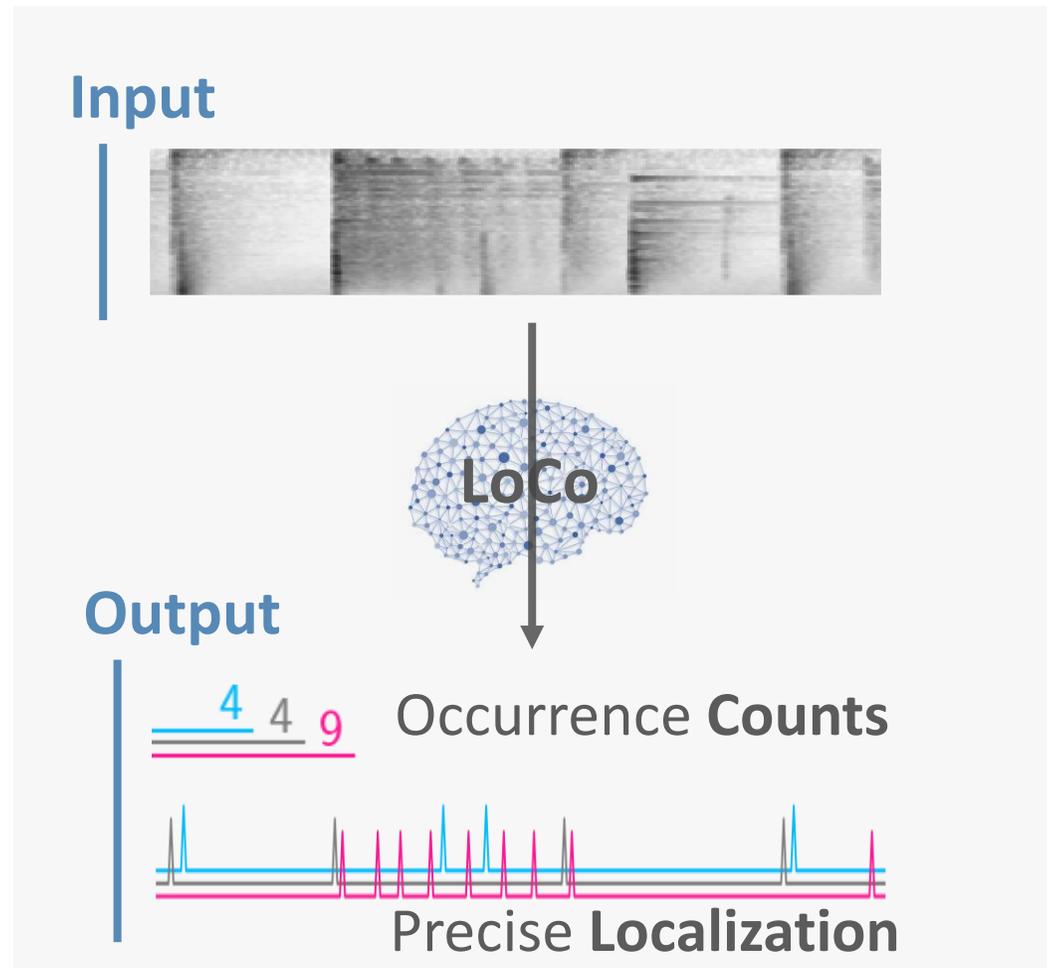
Inference

OBJECTIVE

Weakening the annotation requirement



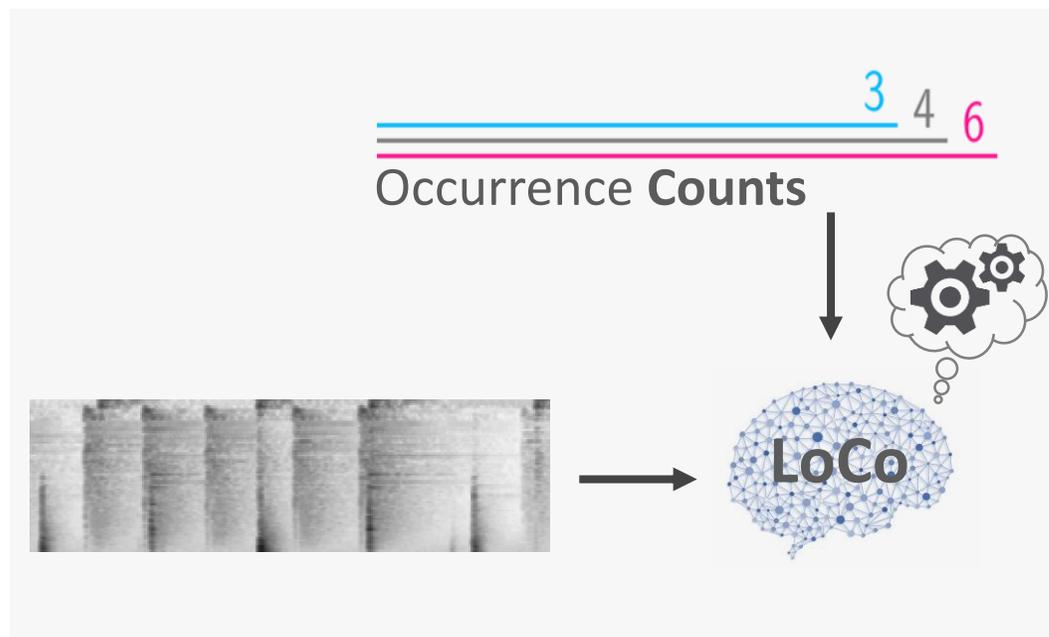
Training



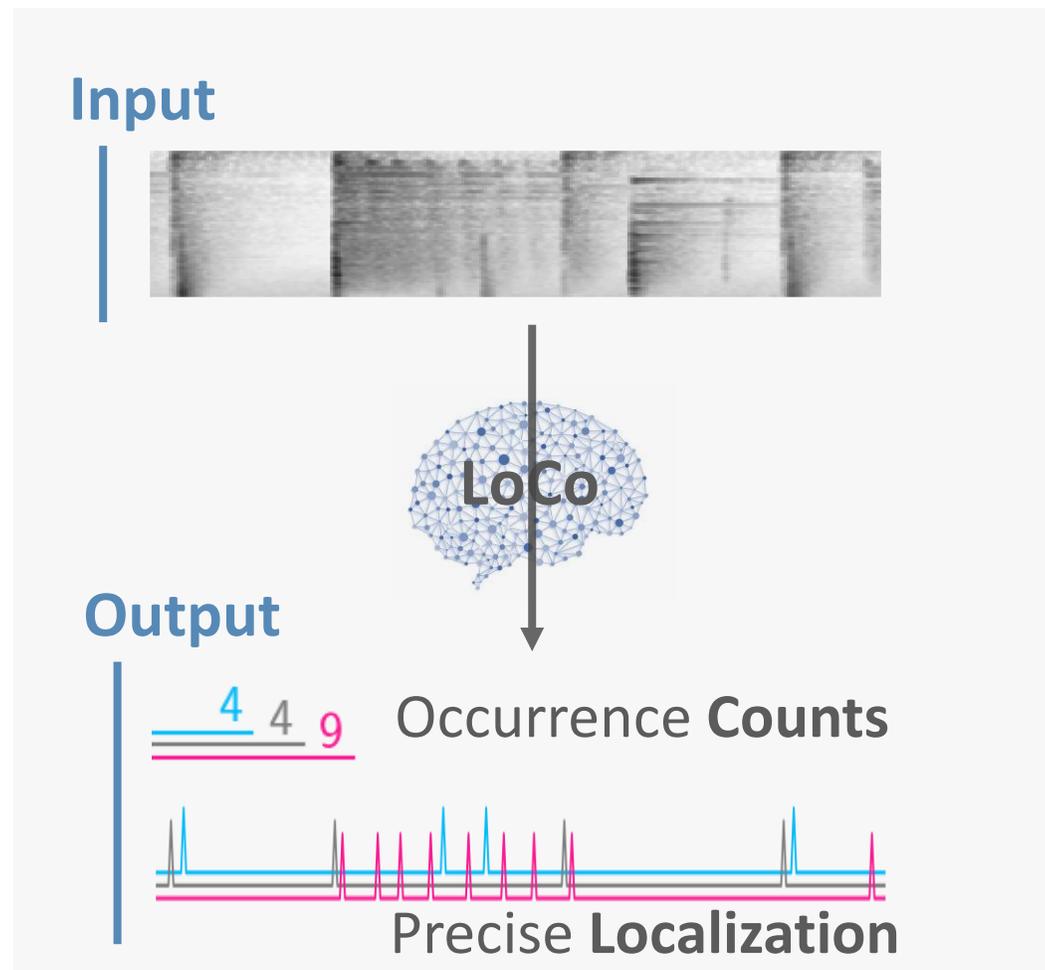
Inference

OBJECTIVE

Weakening the annotation requirement



Training



Inference

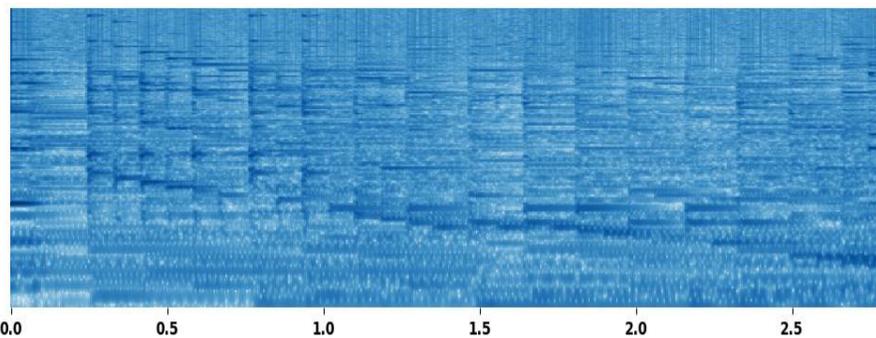
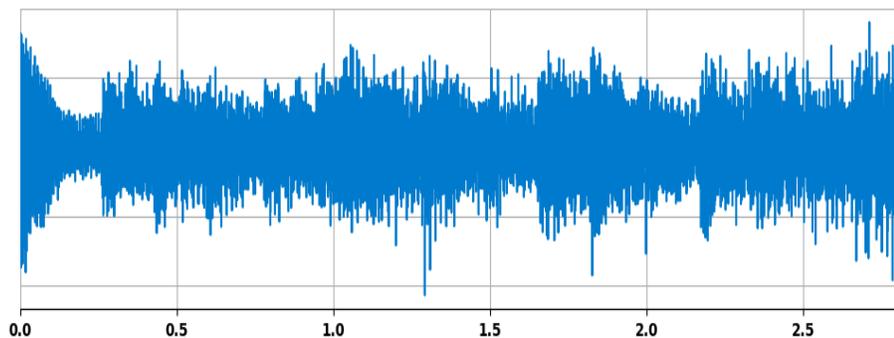
Weakly-Supervised

Is it useful?



OBJECTIVE

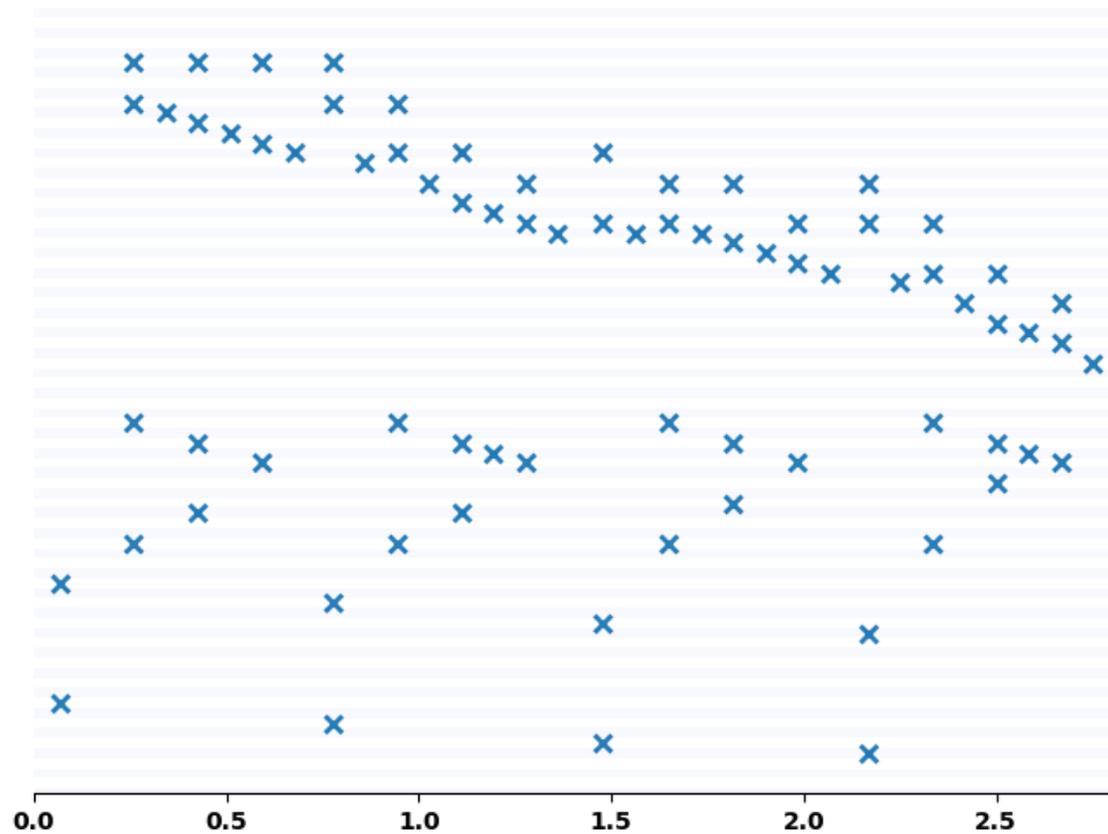
Weakening the annotation requirement



Label Piano Music

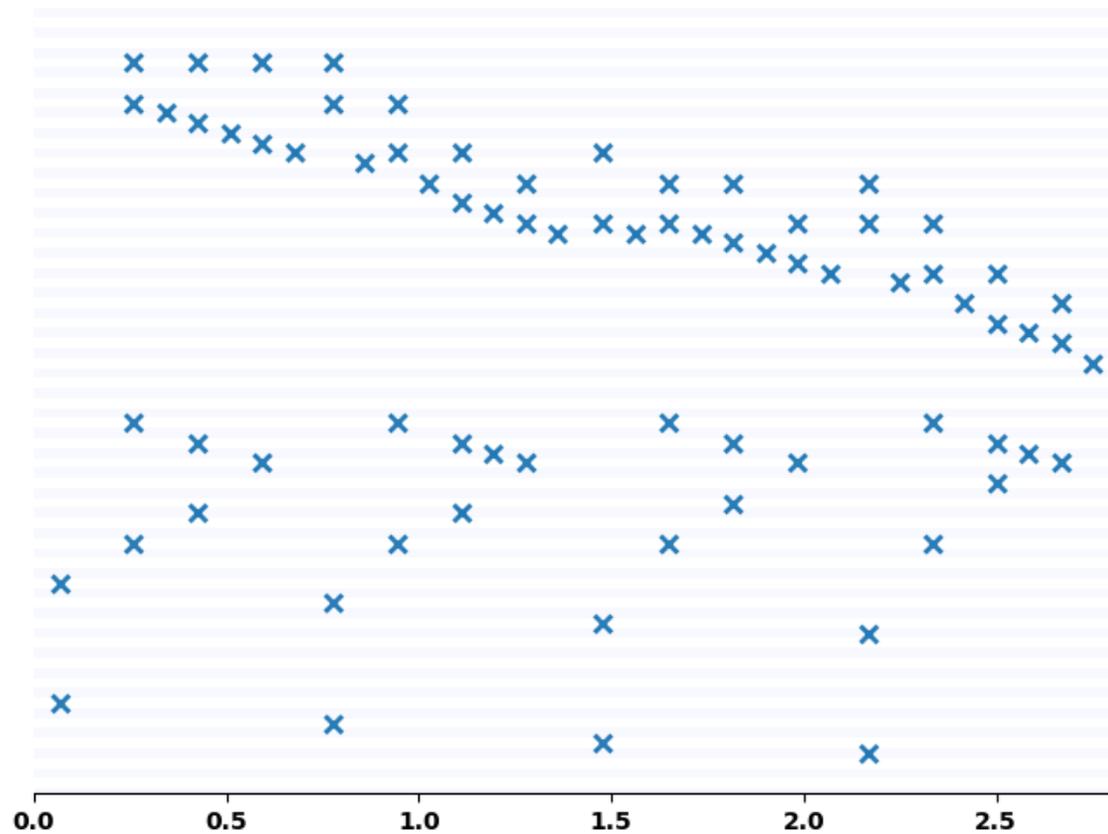
OBJECTIVE

Weakening the annotation requirement



OBJECTIVE

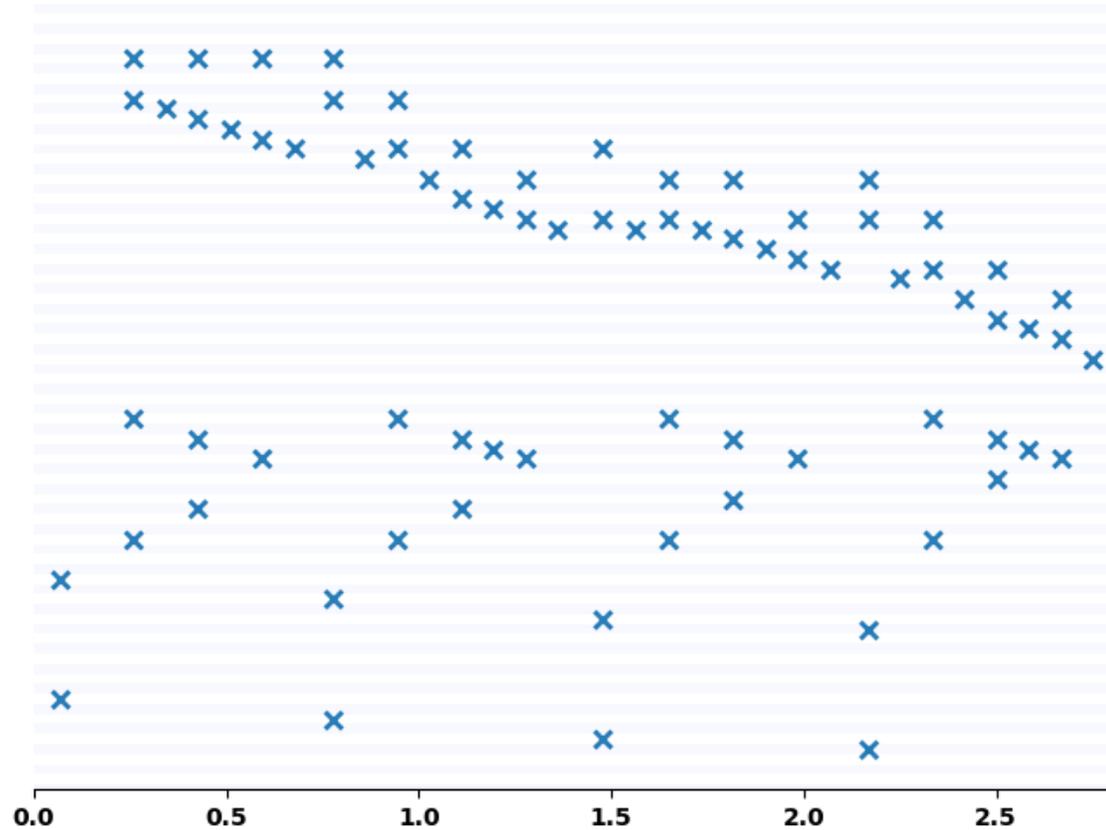
Weakening the annotation requirement



✘ Precise hand-labeling is very **tedious**

OBJECTIVE

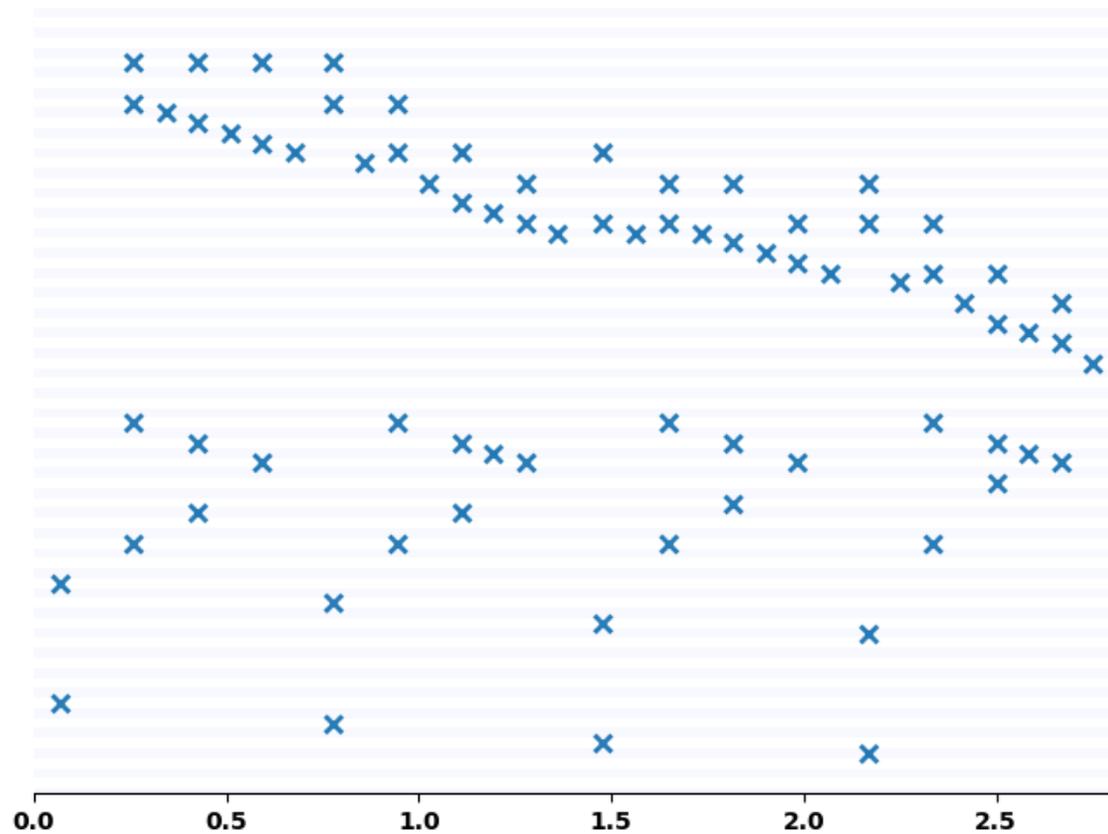
Weakening the annotation requirement



- ✘ Precise hand-labeling is very **tedious**
- ✘ Prone to labeling **inaccuracy**

OBJECTIVE

Weakening the annotation requirement



—————→
Proposed Approach

4
1³
1¹
1⁴
1⁴
5
1¹
1¹
3⁶
1¹
1¹
1³
1²
1¹
1¹
1
4
4²
1
1²
4
1
1
1¹
1
1
1¹

OBJECTIVE

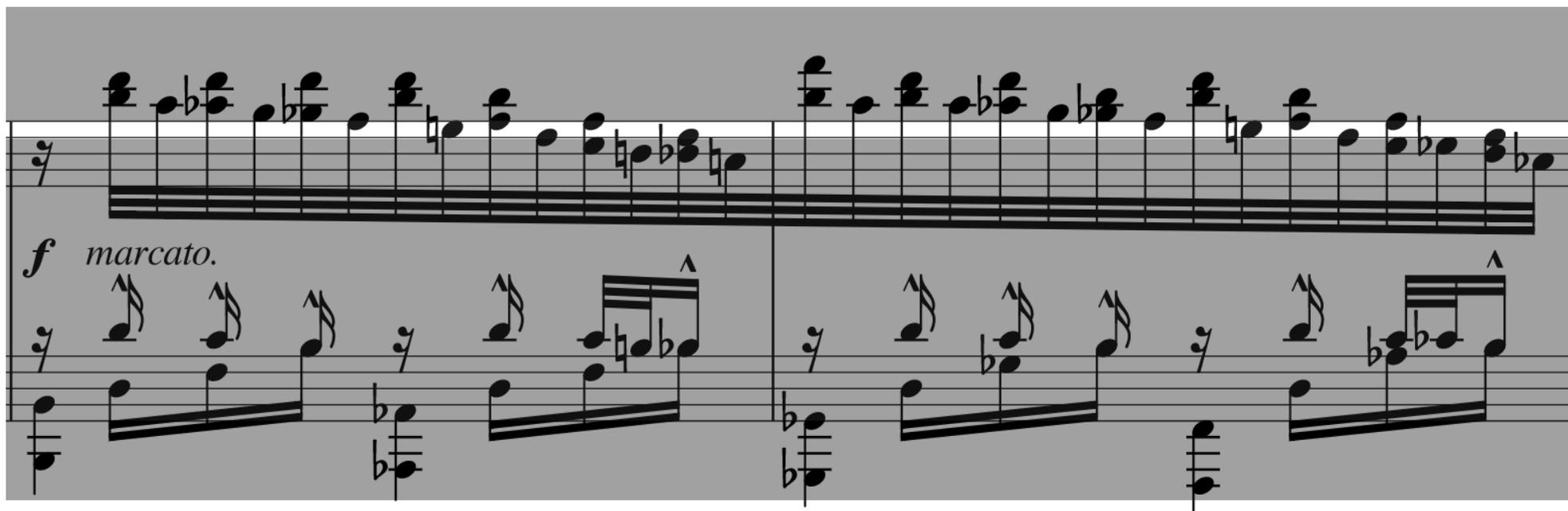
Weakening the annotation requirement

f marcato.

How many notes per pitch?

OBJECTIVE

Weakening the annotation requirement

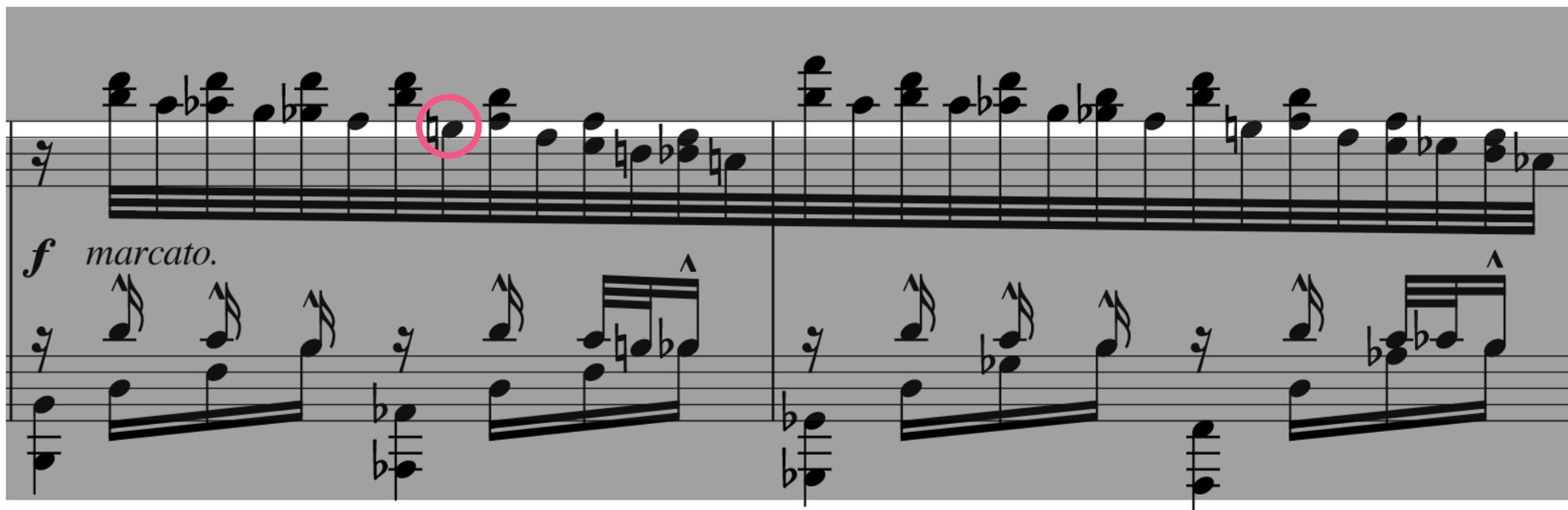


The image shows a musical score for piano, consisting of three staves. The top staff is a single melodic line with many notes, some of which are beamed together. The middle staff is marked *f marcato.* and contains several chords and single notes. The bottom staff contains a bass line with many notes, some of which are beamed together. The score is set against a grey background.

How many notes per pitch?

OBJECTIVE

Weakening the annotation requirement

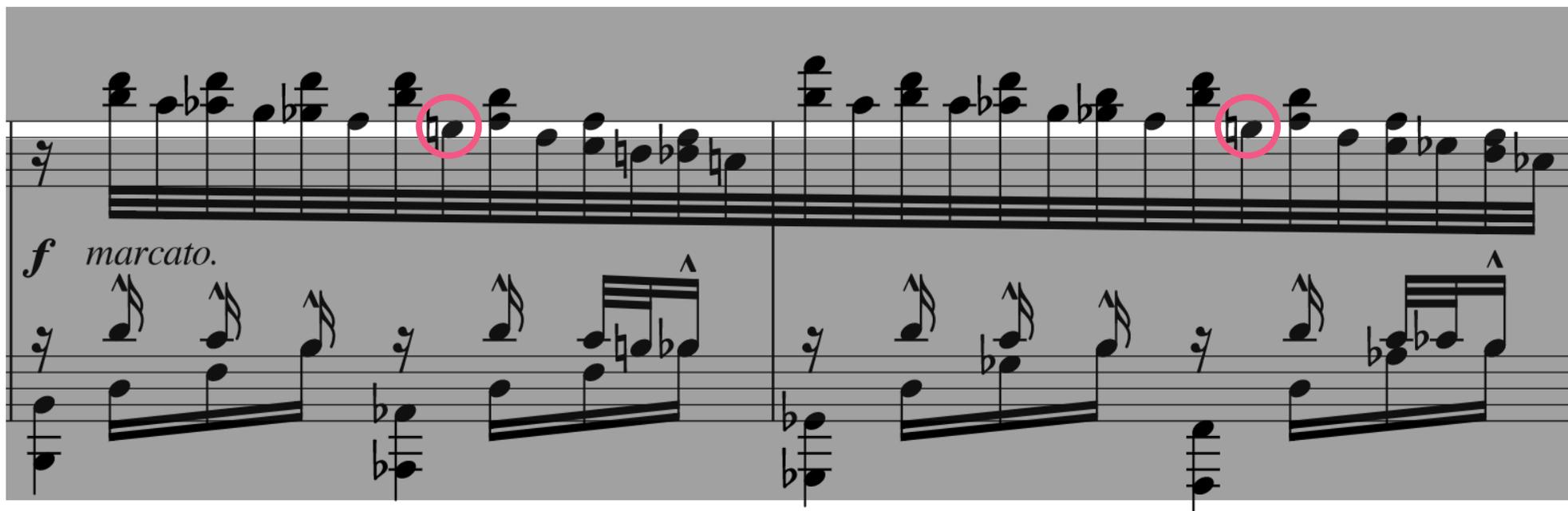


A musical score consisting of three staves. The top staff is a single melodic line with various notes and rests. A red circle highlights a specific note on this staff. The middle and bottom staves are a piano accompaniment, with the middle staff starting with the instruction *f marcato.* The piano part features chords and moving lines in both hands.

How many notes per pitch?

OBJECTIVE

Weakening the annotation requirement

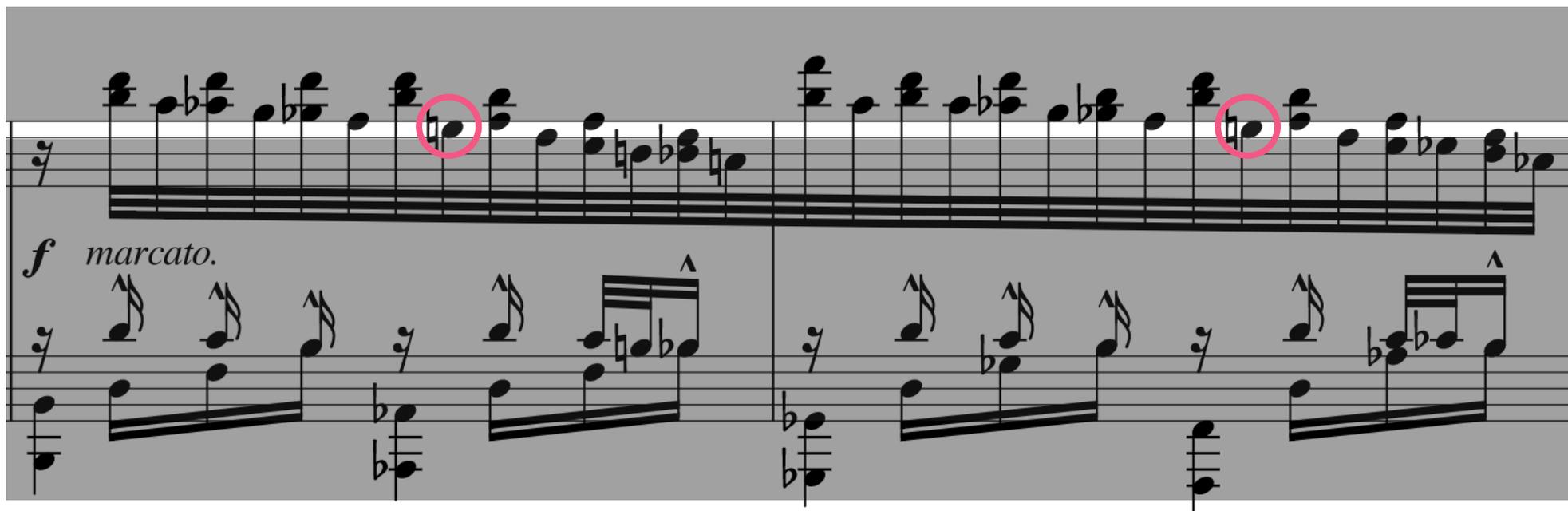


A musical score consisting of three staves. The top staff is a single melodic line with a treble clef and a 7/8 time signature. Two notes in this staff are circled in red. The middle and bottom staves are a piano accompaniment with a bass clef and a 7/8 time signature. The middle staff begins with the dynamic marking *f marcato.* The score is divided into two measures by a vertical bar line.

How many notes per pitch?

OBJECTIVE

Weakening the annotation requirement



A musical score consisting of three staves. The top staff is a single melodic line with several notes circled in red. The middle staff is a piano accompaniment starting with the instruction *f marcato.* The bottom staff is a bass line. A red circle containing the number '2' is located on the right side of the top staff.

How many notes per pitch?

The Model





MODEL

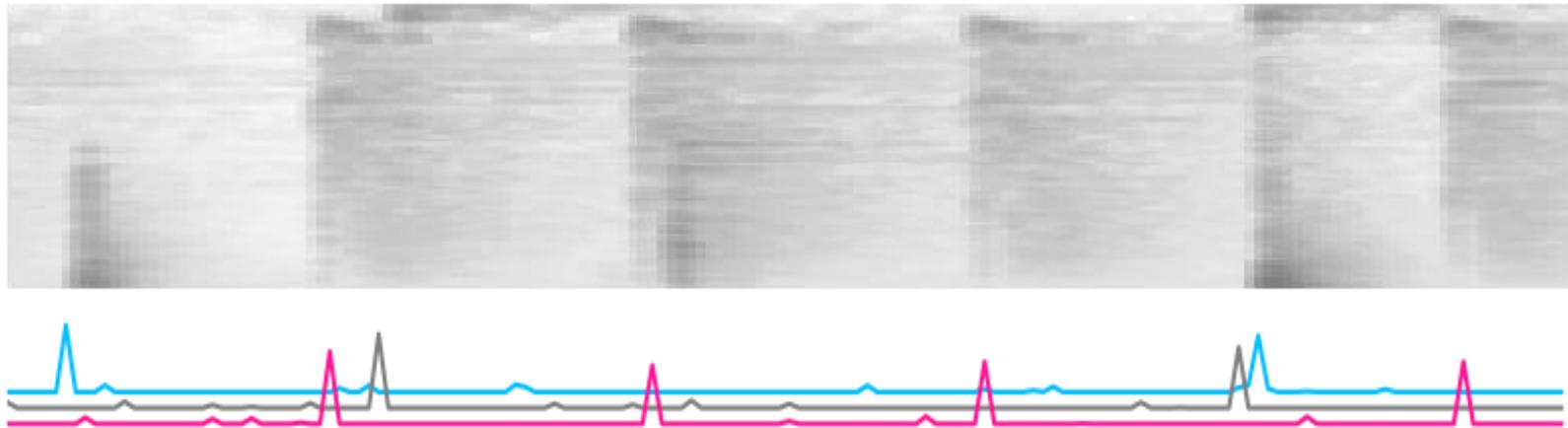
Main Idea

Unlike existing methods, in which localization is explicitly achieved by design, our model learns localization **implicitly** as a byproduct of learning to count instances.

MODEL

Counting Occurrences

$$p_i(t) = f\left(\left(\mathbf{x}_i(n)\right)_{n=1}^t\right) \quad \textit{Probability of Event occurrence}$$





MODEL

Counting Occurrences

$$p_i(t) = f \left(\left(\mathbf{x}_i(n) \right)_{n=1}^t \right)$$

$E_i(t) = \mathfrak{B}(p_i(t))$, ind. Bernoulli *Event occurrence*

MODEL

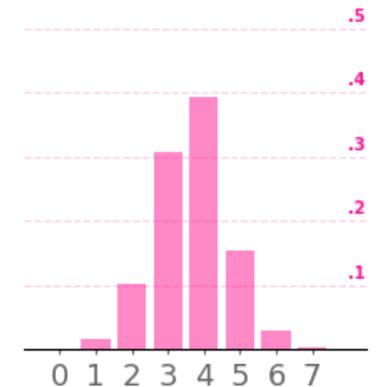
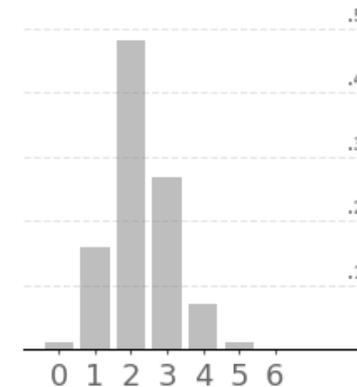
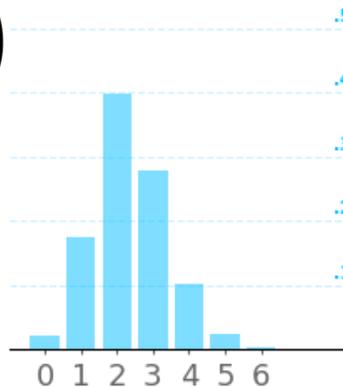
Counting Occurrences

$$p_i(t) = f\left(\left(\mathbf{x}_i(n)\right)_{n=1}^t\right)$$

$$E_i(t) = \mathfrak{B}(p_i(t)), \text{ ind. Bernoulli}$$

$$Y_i = \sum_t E_i(t)$$

Occurrence Count



MODEL

Counting Occurrences

$$p_i(t) = f \left(\overset{\text{Input Data}}{\left(\mathbf{x}_i(n) \right)_{n=1}^t} \right)$$

$$E_i(t) = \mathfrak{B}(p_i(t)), \text{ ind. Bernoulli}$$

$$Y_i = \sum_t E_i(t)$$

MODEL

Counting Occurrences

$$p_i(t) = f \left(\left(\mathbf{x}_i(n) \right)_{n=1}^t \right)$$

$$E_i(t) = \mathfrak{B}(p_i(t)), \text{ ind. Bernoulli}$$

$$Y_i = \sum_t E_i(t)$$

Occurrence
Count

MODEL

Counting Occurrences

$$p_i(t) = f \left(\left(\mathbf{x}_i(n) \right)_{n=1}^t \right)$$

$$E_i(t) = \mathfrak{B}(p_i(t)), \text{ ind. Bernoulli}$$

$$Y_i = \sum_t E_i(t)$$

*Occurrence
Count*

MODEL

Counting Occurrences

Estimated through RNN (e.g. LSTM)
Input Data

$$p_i(t) = f \left(\left(\mathbf{x}_i(n) \right)_{n=1}^t \right)$$

$$E_i(t) = \mathfrak{B}(p_i(t)), \text{ ind. Bernoulli}$$

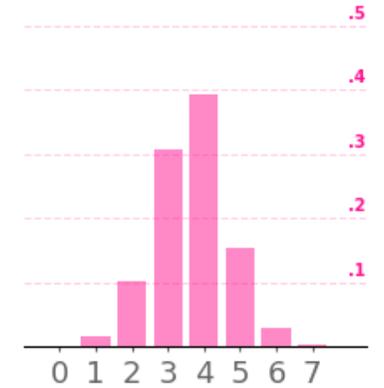
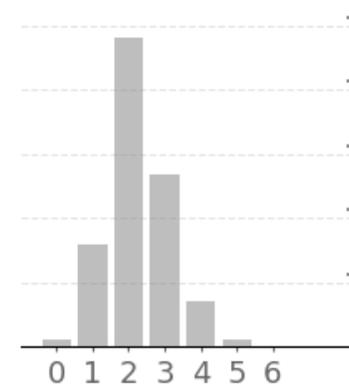
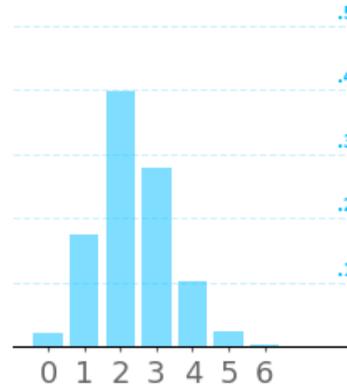
$$Y_i = \sum_t E_i(t)$$

Occurrence
Count

MODEL Loss

$$Y_i = \sum_t E_i(t)$$

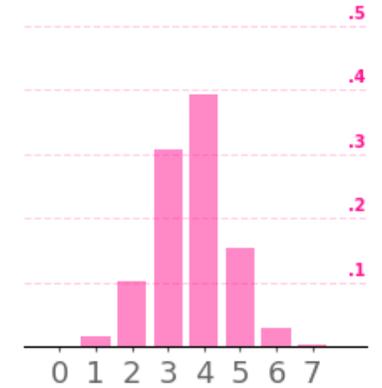
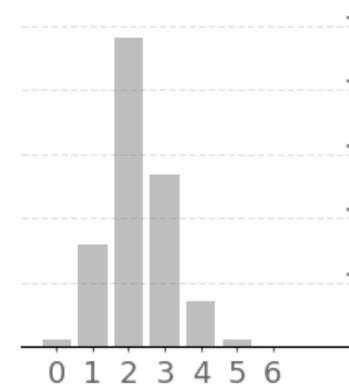
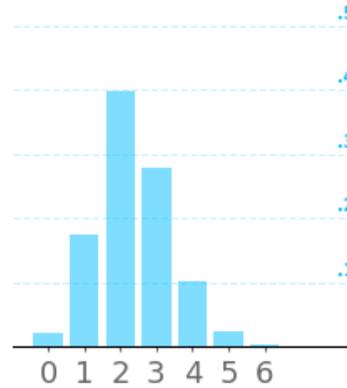
*Occurrence
Count*



MODEL Loss

$$Y_i = \sum_t E_i(t)$$

*Occurrence
Count*

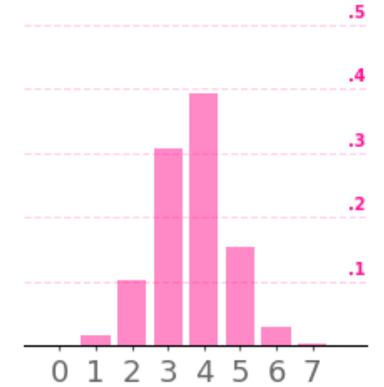
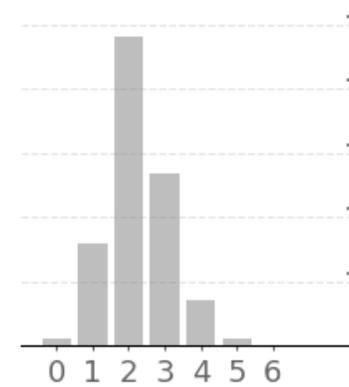
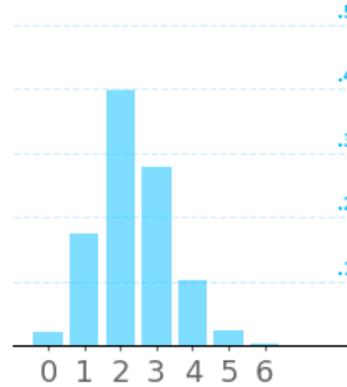


Compare them to true observed counts.

MODEL Loss

$$Y_i = \sum_t E_i(t)$$

*Occurrence
Count*



Compare them to true observed counts.

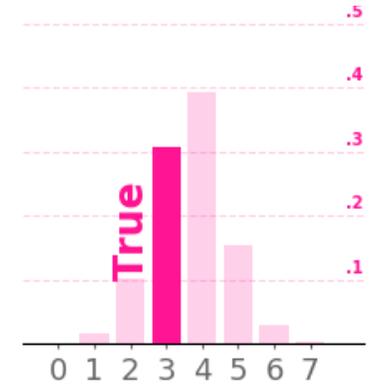
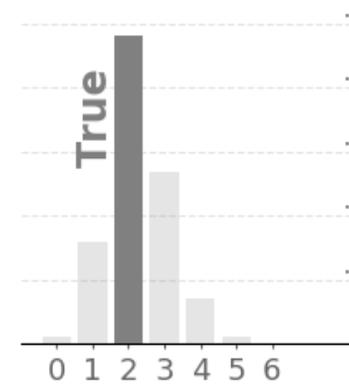
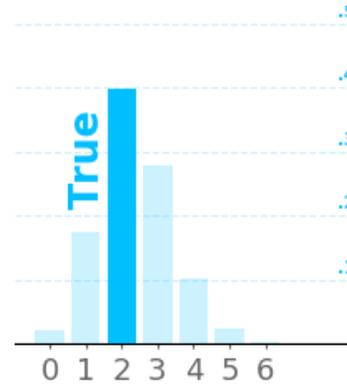
$$L(\theta) = - \sum \log (\text{Pr} (Y_{i,\theta} = y_i | \mathbf{X}_i))$$

Observed Count

MODEL Loss

$$Y_i = \sum_t E_i(t)$$

*Occurrence
Count*



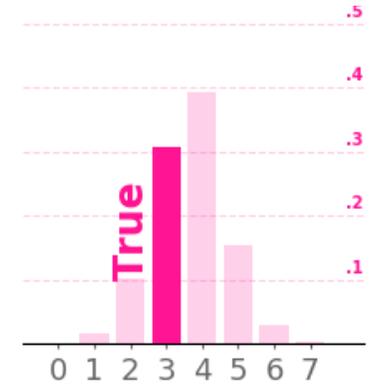
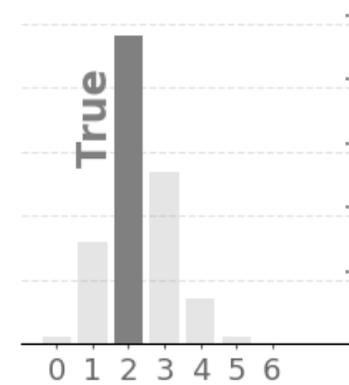
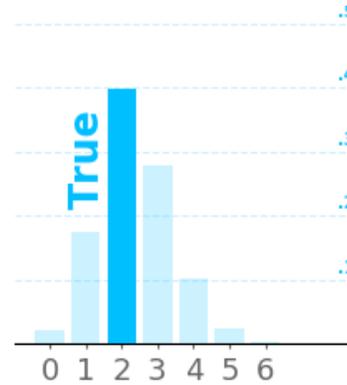
Compare them to true observed counts.

$$L(\theta) = - \sum \log (\text{Pr} (Y_{i,\theta} = y_i | \mathbf{X}_i))$$

MODEL Loss

$$Y_i = \sum_t E_i(t)$$

*Occurrence
Count*



Compare them to true observed counts.

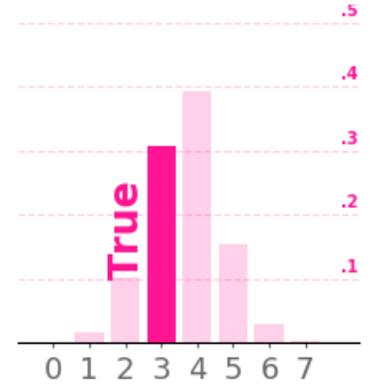
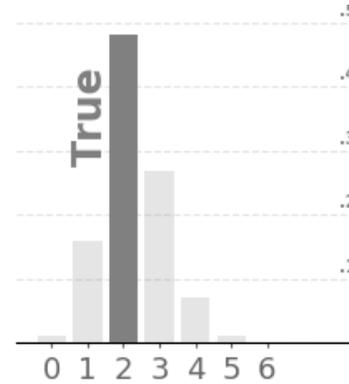
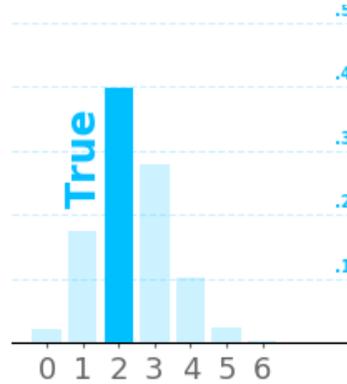
$$L(\theta) = - \sum \log (\text{Pr} (Y_{i,\theta} = y_i | \mathbf{X}_i))$$

$$L(\theta) = -\log \left(\begin{array}{c} \text{blue bar} \\ \text{ } \end{array} \right) - \log \left(\begin{array}{c} \text{grey bar} \\ \text{ } \end{array} \right) - \log \left(\begin{array}{c} \text{pink bar} \\ \text{ } \end{array} \right).$$

MODEL Loss

$$Y_i = \sum_t E_i(t)$$

*Occurrence
Count*



Compare them to true observed counts.

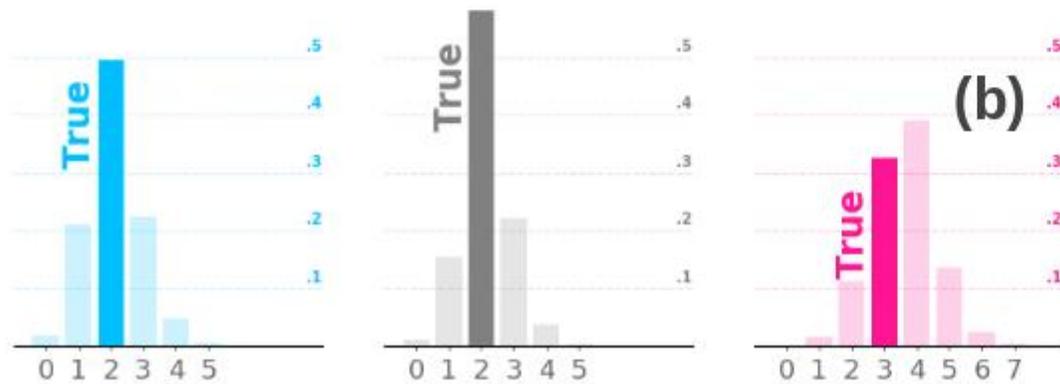
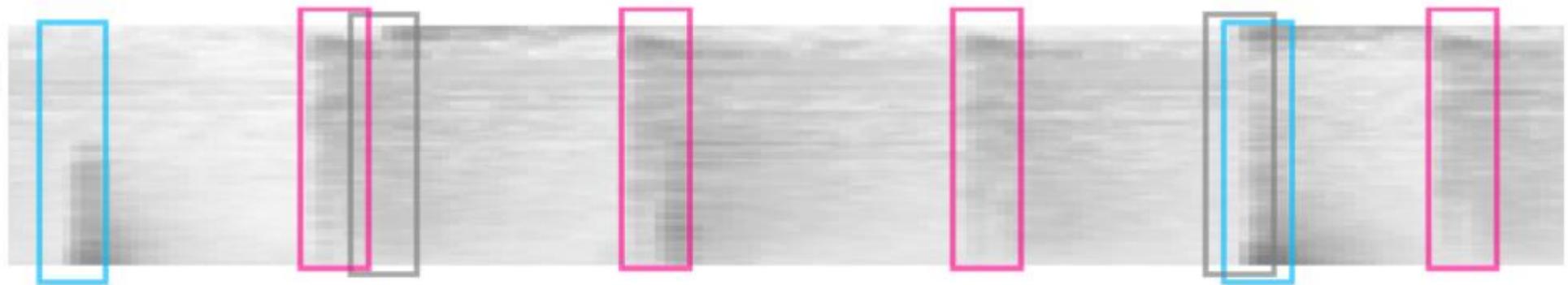
$$L(\theta) = - \sum \log (\text{Pr} (Y_{i,\theta} = y_i | \mathbf{X}_i))$$

Optimized with standard
backpropagation

$$L(\theta) = -\log \left(\begin{matrix} \text{blue bar} \\ \text{ } \end{matrix} \right) - \log \left(\begin{matrix} \text{grey bar} \\ \text{ } \end{matrix} \right) - \log \left(\begin{matrix} \text{pink bar} \\ \text{ } \end{matrix} \right).$$

MODEL

Full Pipeline



$$L(\theta) = -\log \left(\begin{array}{c} \text{blue bar} \\ \text{grey bar} \\ \text{pink bar} \end{array} \right) - \log \left(\begin{array}{c} \text{grey bar} \\ \text{pink bar} \end{array} \right) - \log \left(\begin{array}{c} \text{pink bar} \end{array} \right)$$

(c)

Why does it work?





MODEL

Poisson Binomial Counts

$$\Pr(Y_{i,\theta} = k \mid \mathbf{X}_i) = \sum_{A \in F_k} \prod_{l \in A} \hat{p}_{i,\theta}(l) \prod_{j \in A^c} (1 - \hat{p}_{i,\theta}(j)),$$

Y follows a Poisson-binomial distribution



MODEL Recursion

$$\Upsilon_i(k, t) := \Pr(Y_{i,\theta}(t) = k)$$



*Bin k of count
distribution at time t*

MODEL

Recursion

$$\Upsilon_i(k, t) := \Pr(Y_{i,\theta}(t) = k)$$



*Bin k of count
distribution at time t*

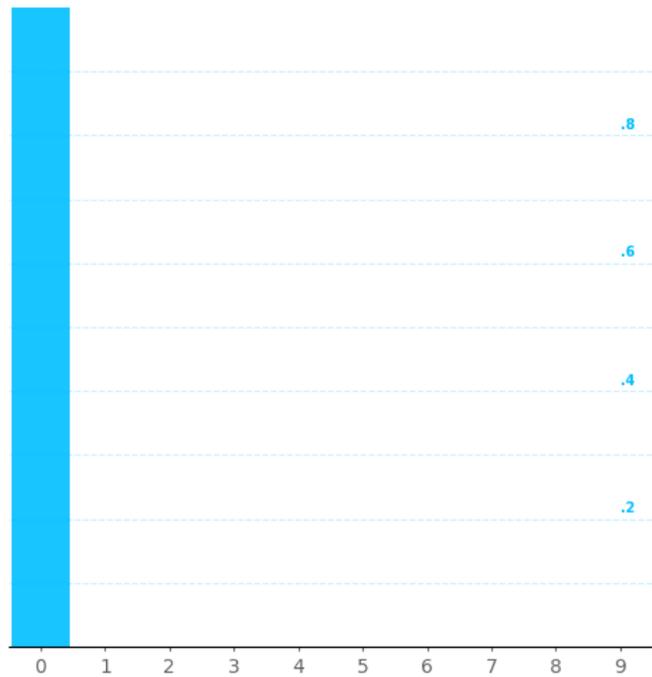
Property 2 (Recursion on k, t)

$$\Upsilon_i(k, t) = \begin{cases} (1 - p_i(t)) \Upsilon_i(k, t-1) & k=0 \\ (1 - p_i(t)) \Upsilon_i(k, t-1) + p_i(t) \Upsilon_i(k-1, t-1) & k>0 \end{cases} \quad (9)$$

where $\Upsilon_i(k, 0) = \mathbb{1}_{k=0}$.

MODEL Recursion

$t = 0$

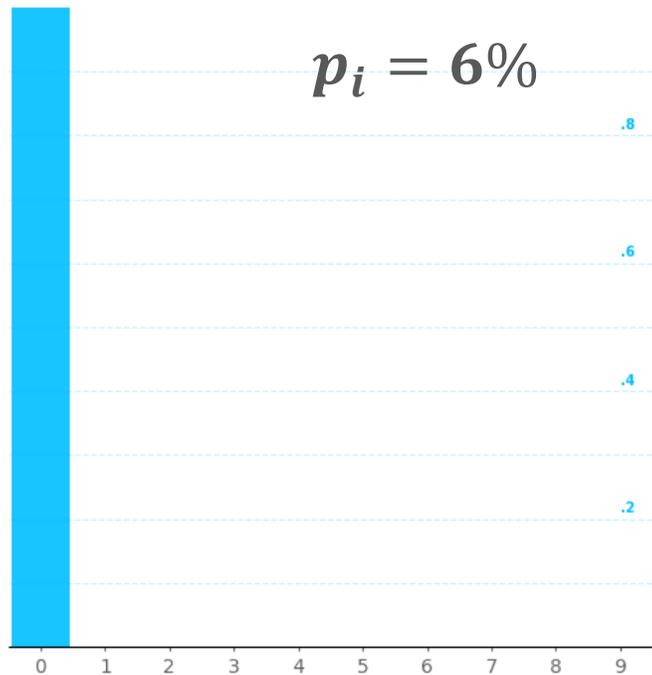


Property 2 (Recursion on k, t)

$$\Upsilon_i(k, t) = \begin{cases} (1-p_i(t))\Upsilon_i(k, t-1) & k=0 \\ (1-p_i(t))\Upsilon_i(k, t-1) + p_i(t)\Upsilon_i(k-1, t-1) & k>0 \end{cases} \quad (9)$$

where $\Upsilon_i(k, 0) = \mathbb{1}_{k=0}$.

MODEL Recursion

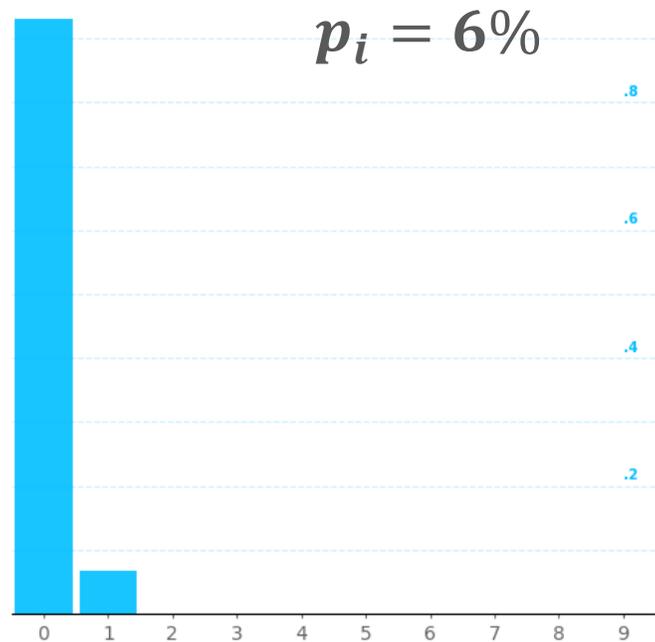


Property 2 (Recursion on k, t)

$$\Upsilon_i(k, t) = \begin{cases} (1 - p_i(t)) \Upsilon_i(k, t-1) & k=0 \\ (1 - p_i(t)) \Upsilon_i(k, t-1) + p_i(t) \Upsilon_i(k-1, t-1) & k > 0 \end{cases} \quad (9)$$

where $\Upsilon_i(k, 0) = \mathbb{1}_{k=0}$.

MODEL Recursion

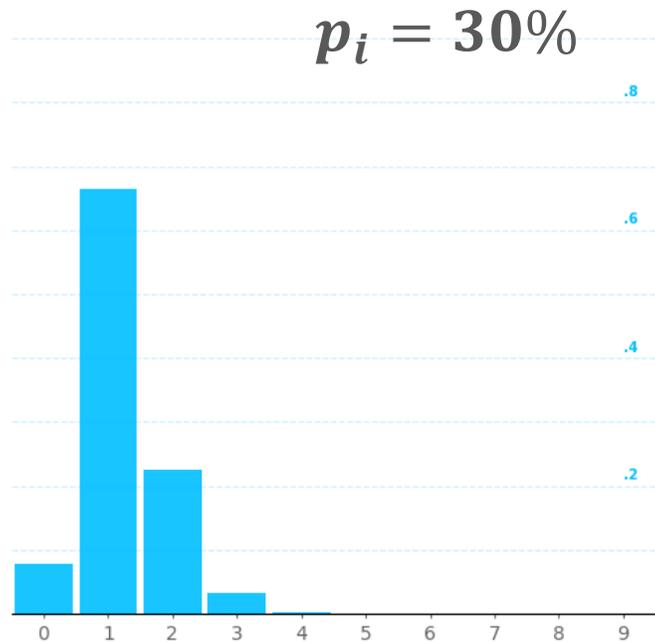


Property 2 (Recursion on k, t)

$$\Upsilon_i(k, t) = \begin{cases} (1 - p_i(t)) \Upsilon_i(k, t-1) & k=0 \\ (1 - p_i(t)) \Upsilon_i(k, t-1) + p_i(t) \Upsilon_i(k-1, t-1) & k > 0 \end{cases} \quad (9)$$

where $\Upsilon_i(k, 0) = \mathbb{1}_{k=0}$.

MODEL Recursion

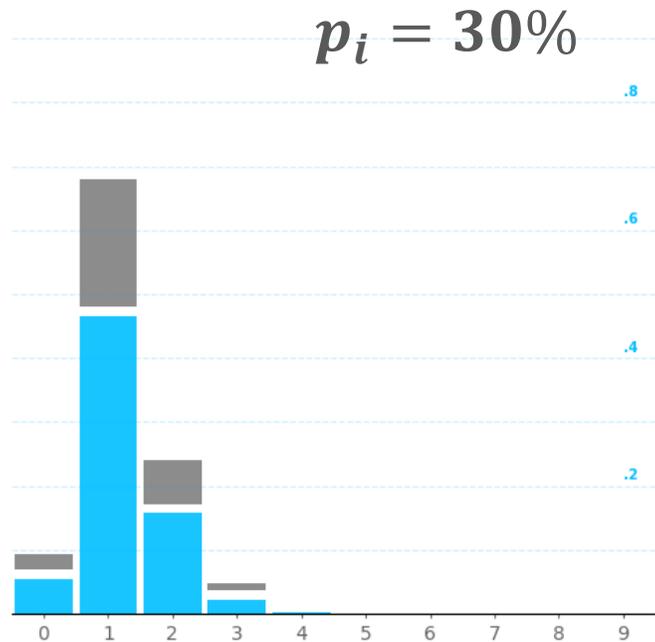


Property 2 (Recursion on k, t)

$$\Upsilon_i(k, t) = \begin{cases} (1-p_i(t))\Upsilon_i(k, t-1) & k=0 \\ (1-p_i(t))\Upsilon_i(k, t-1) + p_i(t)\Upsilon_i(k-1, t-1) & k>0 \end{cases} \quad (9)$$

where $\Upsilon_i(k, 0) = \mathbb{1}_{k=0}$.

MODEL Recursion

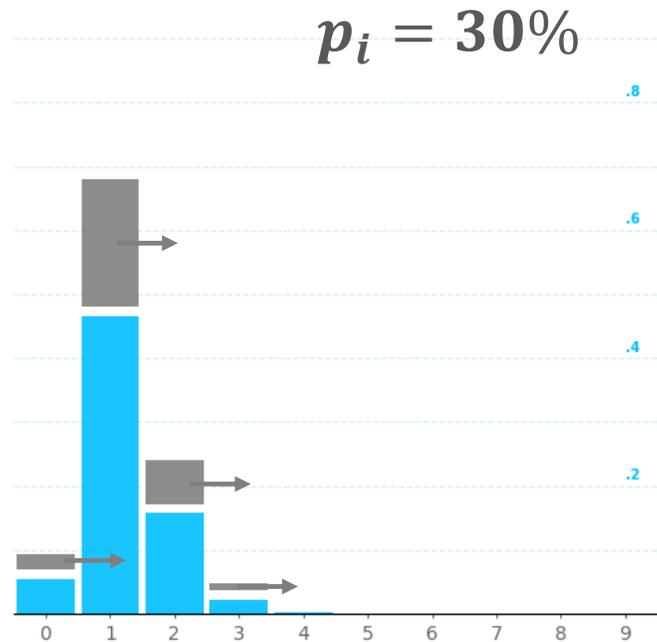


Property 2 (Recursion on k, t)

$$\Upsilon_i(k, t) = \begin{cases} (1-p_i(t))\Upsilon_i(k, t-1) & k=0 \\ (1-p_i(t))\Upsilon_i(k, t-1) + p_i(t)\Upsilon_i(k-1, t-1) & k>0 \end{cases} \quad (9)$$

where $\Upsilon_i(k, 0) = \mathbb{1}_{k=0}$.

MODEL Recursion

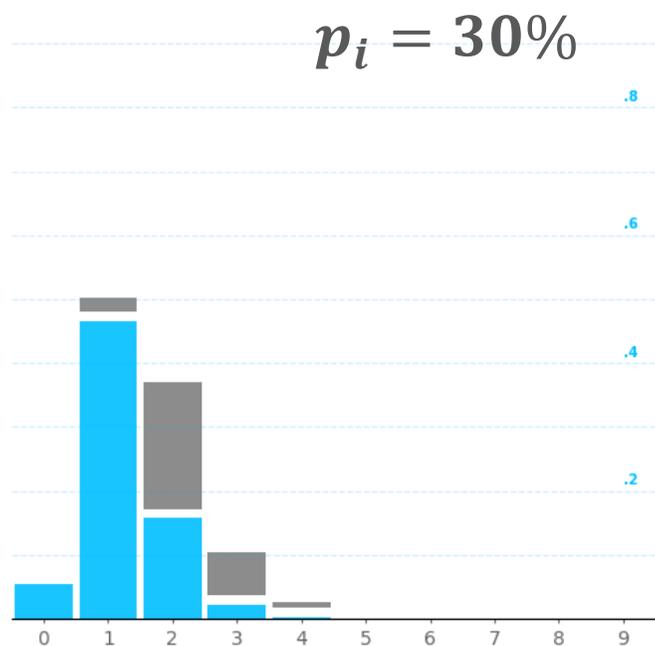


Property 2 (Recursion on k, t)

$$\Upsilon_i(k, t) = \begin{cases} (1-p_i(t))\Upsilon_i(k, t-1) & k=0 \\ (1-p_i(t))\Upsilon_i(k, t-1) + p_i(t)\Upsilon_i(k-1, t-1) & k>0 \end{cases} \quad (9)$$

where $\Upsilon_i(k, 0) = \mathbb{1}_{k=0}$.

MODEL Recursion

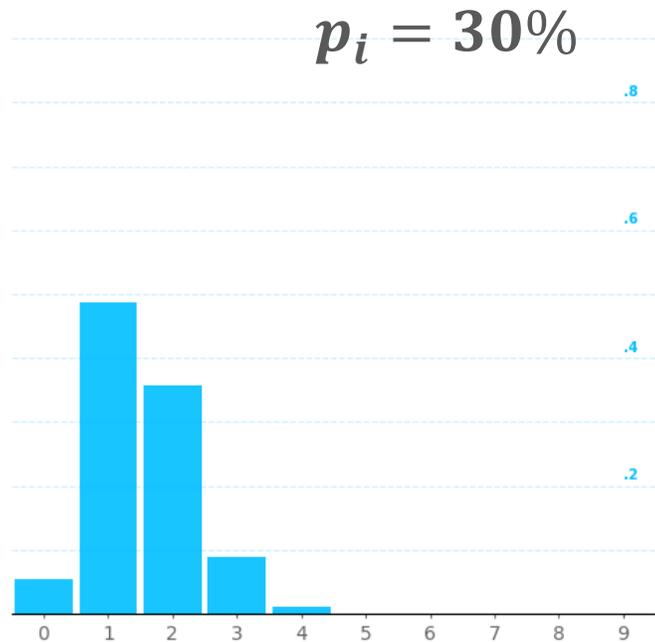


Property 2 (Recursion on k, t)

$$\Upsilon_i(k, t) = \begin{cases} (1 - p_i(t)) \Upsilon_i(k, t-1) & k=0 \\ (1 - p_i(t)) \Upsilon_i(k, t-1) + p_i(t) \Upsilon_i(k-1, t-1) & k > 0 \end{cases} \quad (9)$$

where $\Upsilon_i(k, 0) = \mathbb{1}_{k=0}$.

MODEL Recursion



Property 2 (Recursion on k, t)

$$\Upsilon_i(k, t) = \begin{cases} (1-p_i(t))\Upsilon_i(k, t-1) & k=0 \\ (1-p_i(t))\Upsilon_i(k, t-1) + p_i(t)\Upsilon_i(k-1, t-1) & k>0 \end{cases} \quad (9)$$

where $\Upsilon_i(k, 0) = \mathbb{1}_{k=0}$.



MODEL

No early triggering

Property 1 (Mass shift irreversibility)

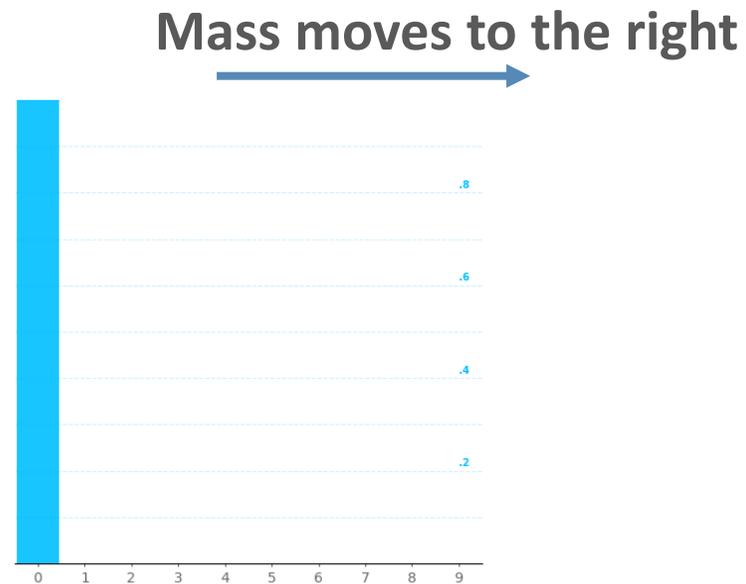
$(Y_{i,\theta}(t))_{t=1}^{T_i}$ is monotonically increasing.

MODEL

No early triggering

Property 1 (Mass shift irreversibility)

$(Y_{i,\theta}(t))_{t=1}^{T_i}$ is monotonically increasing.



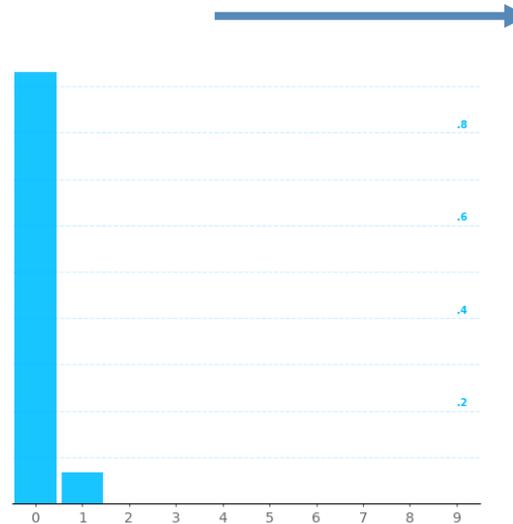
MODEL

No early triggering

Property 1 (Mass shift irreversibility)

$(Y_{i,\theta}(t))_{t=1}^{T_i}$ is monotonically increasing.

Mass moves to the right



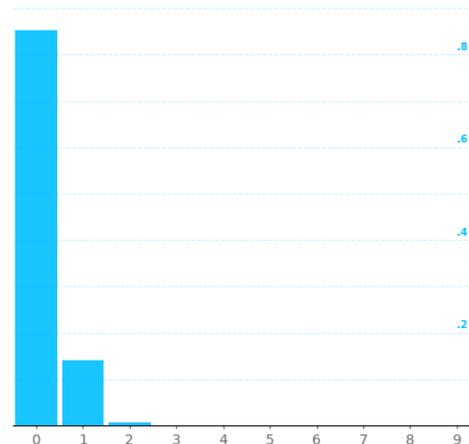
MODEL

No early triggering

Property 1 (Mass shift irreversibility)

$(Y_{i,\theta}(t))_{t=1}^{T_i}$ is monotonically increasing.

Mass moves to the right



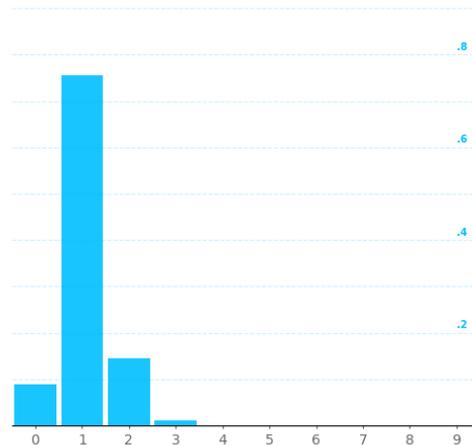
MODEL

No early triggering

Property 1 (Mass shift irreversibility)

$(Y_{i,\theta}(t))_{t=1}^{T_i}$ is monotonically increasing.

Mass moves to the right



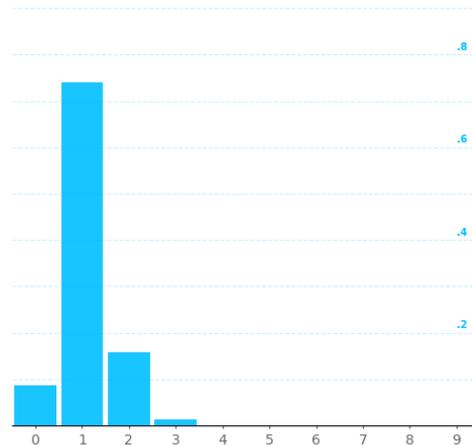
MODEL

No early triggering

Property 1 (Mass shift irreversibility)

$(Y_{i,\theta}(t))_{t=1}^{T_i}$ is monotonically increasing.

Mass moves to the right



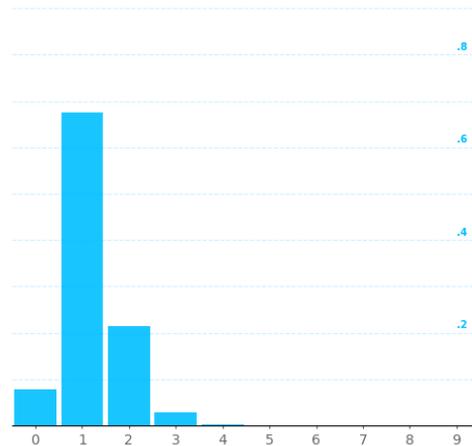
MODEL

No early triggering

Property 1 (Mass shift irreversibility)

$(Y_{i,\theta}(t))_{t=1}^{T_i}$ is monotonically increasing.

Mass moves to the right



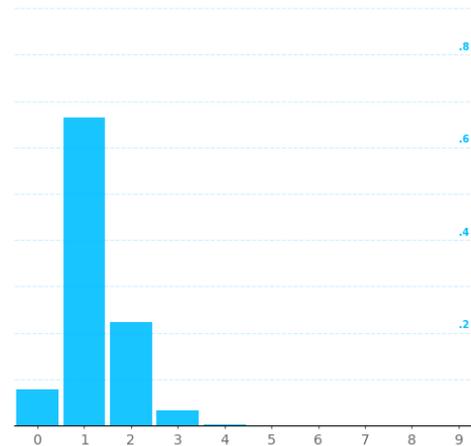
MODEL

No early triggering

Property 1 (Mass shift irreversibility)

$(Y_{i,\theta}(t))_{t=1}^{T_i}$ is monotonically increasing.

Mass moves to the right



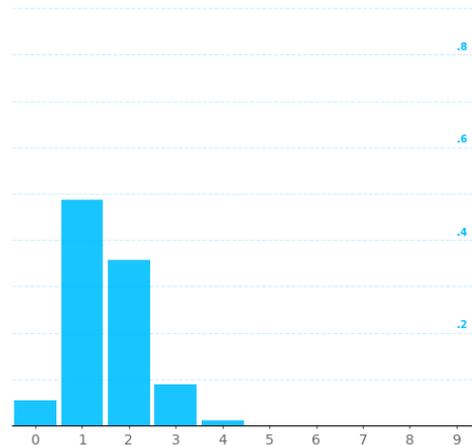
MODEL

No early triggering

Property 1 (Mass shift irreversibility)

$(Y_{i,\theta}(t))_{t=1}^{T_i}$ is monotonically increasing.

Mass moves to the right

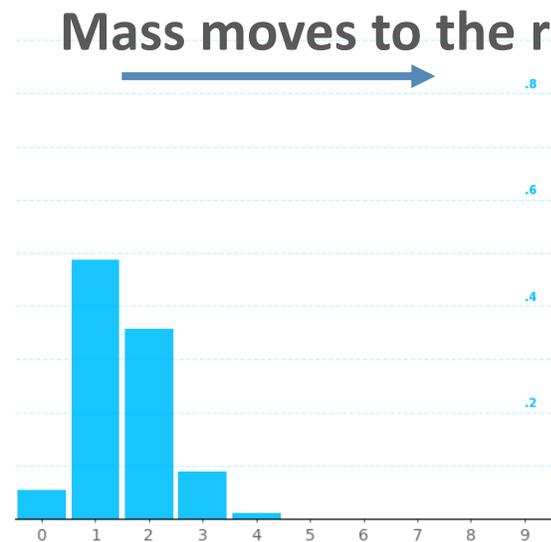


MODEL

No early triggering

Property 1 (Mass shift irreversibility)

$(Y_{i,\theta}(t))_{t=1}^{T_i}$ is monotonically increasing.



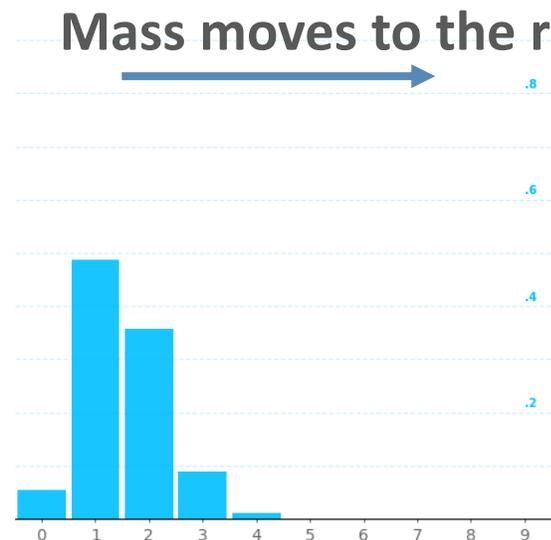
Consequence: Mass shifts are irreversible

MODEL

No early triggering

Property 1 (Mass shift irreversibility)

$(Y_{i,\theta}(t))_{t=1}^{T_i}$ is monotonically increasing.



Consequence: Mass shifts are irreversible

- prevents the model from triggering early
- prevents the model from false alarms



MODEL

Mass Convergence

Lemma 2 (First upper bound)

$$\max_k \Upsilon_i(k, t) \leq \frac{1}{2} + \min_{j \leq t} \left\| \frac{1}{2} - p_i(j) \right\|$$



MODEL

Mass Convergence

Lemma 2 (First upper bound)

$$\max_k \Upsilon_i(k, t) \leq \frac{1}{2} + \min_{j \leq t} \left\| \frac{1}{2} - p_i(j) \right\|$$

$$L(\theta) = - \sum_i \log (\Pr (Y_{i,\theta} = y_i \mid \mathbf{X}_i)) \quad \text{Counting Loss}$$



MODEL

Mass Convergence

Lemma 2 (First upper bound)

$$\max_k \Upsilon_i(k, t) \leq \frac{1}{2} + \min_{j \leq t} \left\| \frac{1}{2} - p_i(j) \right\|$$

$$\begin{aligned} L(\theta) &= - \sum_i \log (\Pr (Y_{i,\theta} = y_i \mid \mathbf{X}_i)) && \text{Counting Loss} \\ &= - \sum_i \log (\Upsilon_i(y_i, T_i)) \end{aligned}$$

MODEL

Mass Convergence

Lemma 2 (First upper bound)

$$\max_k \Upsilon_i(k, t) \leq \frac{1}{2} + \min_{j \leq t} \left\| \frac{1}{2} - p_i(j) \right\|$$

$$L(\theta) = - \sum_i \log (\Pr (Y_{i,\theta} = y_i \mid \mathbf{X}_i)) \quad \text{Counting Loss}$$

$$= - \sum_i \log (\Upsilon_i(y_i, T_i))$$

$$\stackrel{(11)}{\geq} - \sum_i \log \left(\frac{1}{2} + \min_{j \leq t} \left\| \frac{1}{2} - p_i(j) \right\| \right)$$



MODEL

Mass Convergence

Lemma 2 (First upper bound)

$$\max_k \Upsilon_i(k, t) \leq \frac{1}{2} + \min_{j \leq t} \left\| \frac{1}{2} - p_i(j) \right\|$$

Learns to count ↓ $L(\theta) = - \sum_i \log (\Pr (Y_{i,\theta} = y_i \mid \mathbf{X}_i))$ **Counting Loss**

$$= - \sum_i \log (\Upsilon_i(y_i, T_i))$$

$$\stackrel{(11)}{\geq} - \sum_i \log \left(\frac{1}{2} + \min_{j \leq t} \left\| \frac{1}{2} - p_i(j) \right\| \right)$$

MODEL

Mass Convergence

Lemma 2 (First upper bound)

$$\max_k \Upsilon_i(k, t) \leq \frac{1}{2} + \min_{j \leq t} \left\| \frac{1}{2} - p_i(j) \right\|$$

Learns to count ↓ $L(\theta) = - \sum_i \log (\Pr (Y_{i,\theta} = y_i \mid \mathbf{X}_i))$ **Counting Loss**

$$= - \sum_i \log (\Upsilon_i(y_i, T_i))$$

Converge towards 0,1 extremes

$$\stackrel{(11)}{\geq} - \sum_i \log \left(\frac{1}{2} + \min_{j \leq t} \left\| \frac{1}{2} - p_i(j) \right\| \right)$$

MODEL

Mass Convergence

Property 3 (Sparse mass concentration) The inequality derived below reveals that, as the loss decreases, small $p_i(\cdot)$ will quickly converge towards zero.

$$\begin{aligned} \max_k \Upsilon_i(k, t) &\stackrel{(8)}{\leq} \min_{l \leq t} \max_k \Upsilon_i(k, l) \stackrel{\text{ind}}{=} \min_{\sigma, l \leq t} \max_k \Upsilon_{i, \sigma}(k, l) \\ &\stackrel{\text{Le Cam}}{\leq} \min_{\sigma, l \leq t} \max_k \frac{\lambda_{i, \sigma, l}^k e^{-\lambda_{i, \sigma, l}}}{k!} + 2 \sum_{j=1}^l p_{i, \sigma}(j)^2 \\ &\stackrel{\text{def}}{=} \min_{\sigma, l \leq t} \max_k \frac{\left[\sum_{j=1}^l p_{i, \sigma}(j) \right]^k e^{-\left[\sum_{j=1}^l p_{i, \sigma}(j) \right]}}{k!} + 2 \sum_{j=1}^l p_{i, \sigma}(j)^2, \end{aligned} \quad (13)$$

MODEL

Mass Convergence

Property 3 (Sparse mass concentration) The inequality derived below reveals that, as the loss decreases, small $p_i(\cdot)$ will quickly converge towards zero.

$$\begin{aligned} \max_k \Upsilon_i(k, t) &\stackrel{(8)}{\leq} \min_{l \leq t} \max_k \Upsilon_i(k, l) \stackrel{\text{ind}}{=} \min_{\sigma, l \leq t} \max_k \Upsilon_{i, \sigma}(k, l) \\ &\stackrel{\text{Le Cam}}{\leq} \min_{\sigma, l \leq t} \max_k \frac{\lambda_{i, \sigma, l}^k e^{-\lambda_{i, \sigma, l}}}{k!} + 2 \sum_{j=1}^l p_{i, \sigma}(j)^2 \\ &\stackrel{\text{def}}{=} \min_{\sigma, l \leq t} \max_k \frac{\left[\sum_{j=1}^l p_{i, \sigma}(j) \right]^k e^{-\left[\sum_{j=1}^l p_{i, \sigma}(j) \right]}}{k!} + 2 \sum_{j=1}^l p_{i, \sigma}(j)^2, \end{aligned} \quad (13)$$

A detection cannot be split into numerous small $p_i(\cdot)$ contributions



MODEL

Mass Convergence

As the model **learns to count** event occurrences:



MODEL

Mass Convergence

As the model **learns to count** event occurrences:

- $p_i(\cdot)$ converge towards 0,1 extremes



MODEL

Mass Convergence

As the model **learns to count** event occurrences:

- $p_i(\cdot)$ converge towards 0,1 extremes
- A detection cannot be split into numerous small $p_i(\cdot)$ contributions



MODEL

Mass Convergence

As the model **learns to count** event occurrences:

- $p_i(\cdot)$ converge towards 0,1 extremes
- A detection cannot be split into numerous small $p_i(\cdot)$ contributions

A single $p_i(\cdot)$ will contain almost all of them mass for an event.



MODEL Properties

1. Almost **binary** predictions



MODEL Properties

1. Almost **binary** predictions
2. No **early** triggering



MODEL Properties

1. Almost **binary** predictions
2. No **early** triggering
3. No systematic **late bias** ← *Not a theoretical property*



MODEL Properties

1. Almost **binary** predictions
2. No **early** triggering
3. No systematic **late bias** ← *Not a theoretical property*

Achieved through an implementation trick:



MODEL Properties

1. Almost **binary** predictions
2. No **early** triggering
3. No systematic **late bias** ← *Not a theoretical property*

Achieved through an implementation trick:
Feeding sequences of variable length



MODEL Properties

1. Almost **binary** predictions
2. No **early** triggering
3. No systematic **late** bias

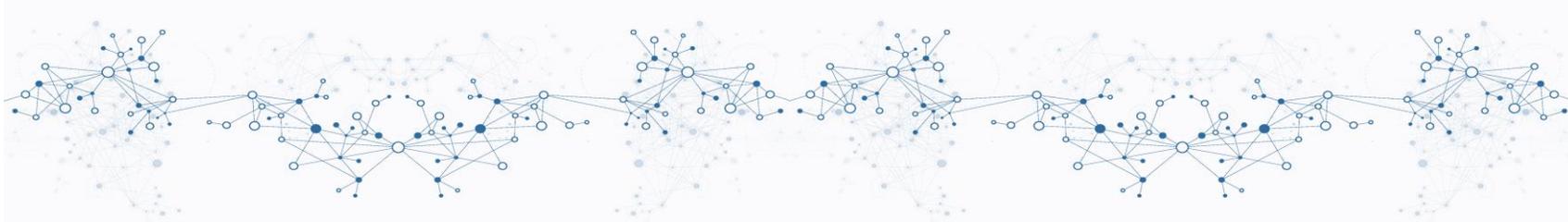


MODEL Properties

1. Almost **binary** predictions
2. No **early** triggering
3. No systematic **late** bias

If the model accurately learns to count occurrences and if the events are detectable, then a coherent localization will emerge naturally.

Experiments





DRUM DETECTION

Experiment Specifications



Detection of three different **drum** types in drum audio extracts



DRUM DETECTION

Experiment Specifications



Detection of three different **drum** types in drum audio extracts

- Tight tolerance of 50ms for a prediction to be correct



DRUM DETECTION

Experiment Specifications



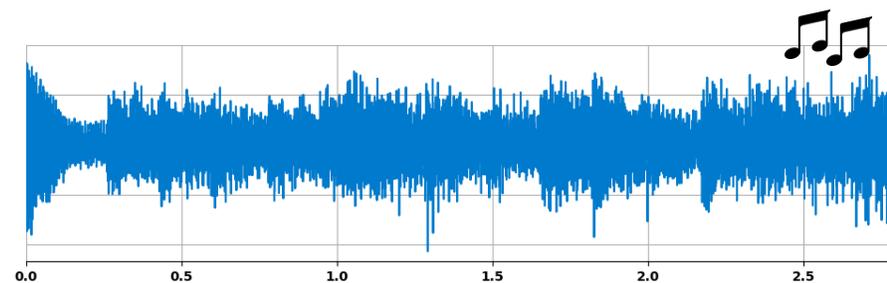
Detection of three different **drum** types in drum audio extracts

- **Tight tolerance** of 50ms for a prediction to be correct
- Comparison with **fully-supervised benchmark** models



DRUM DETECTION

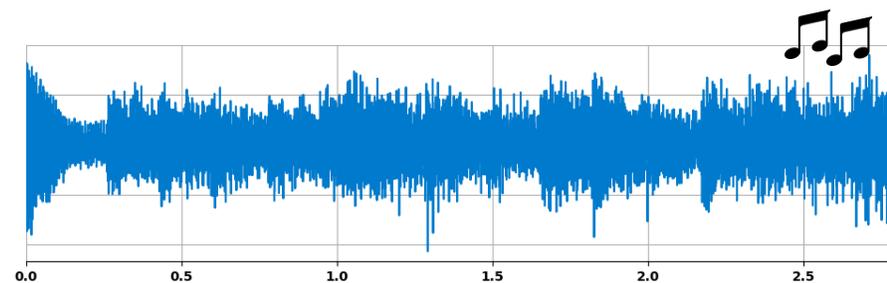
Our approach



Signal

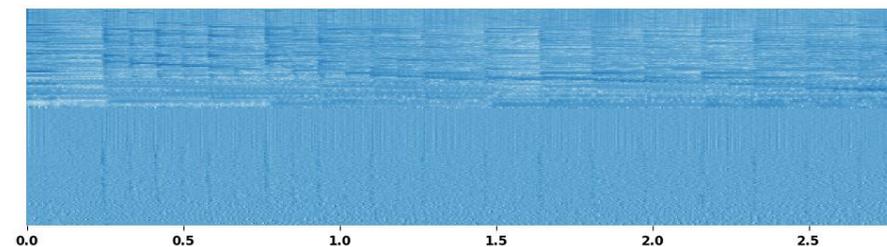
DRUM DETECTION

Our approach



Signal

↓ Fourier

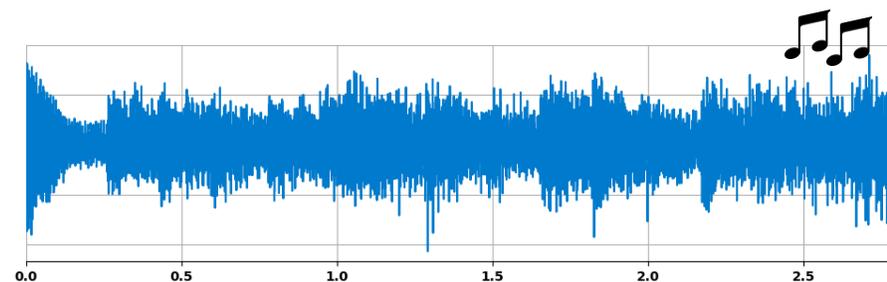


Mel-spectrogram

1st order derivative

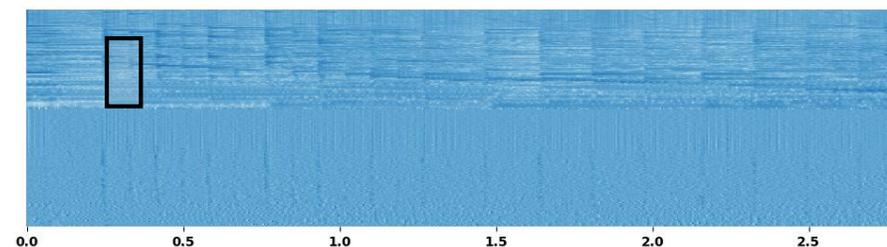
DRUM DETECTION

Our approach



Signal

↓ Fourier



Mel-spectrogram

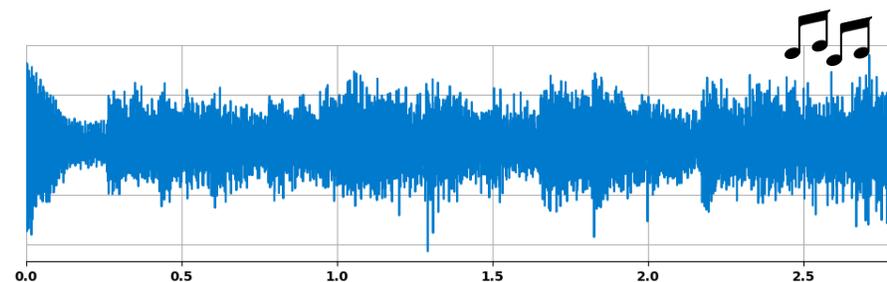
1st order derivative

↓ CNNs

Convolutional Representations

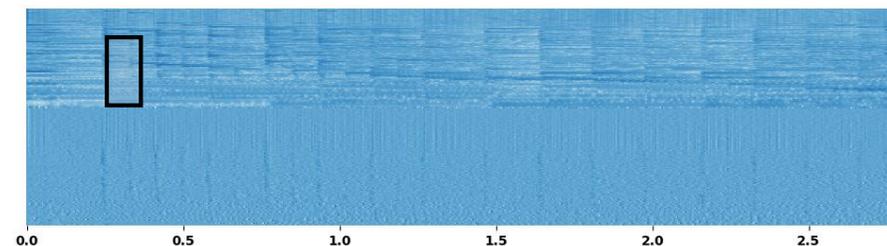
DRUM DETECTION

Our approach



Signal

↓ Fourier



Mel-spectrogram

1st order derivative

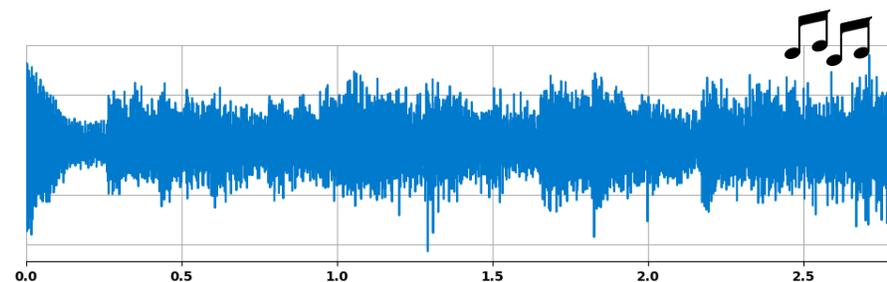
↓ CNNs

Convolutional Representations

↓ LSTM

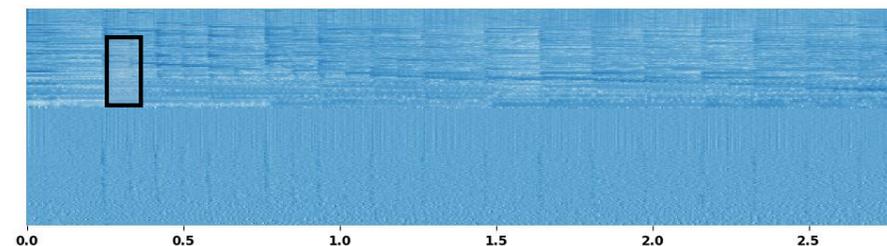
DRUM DETECTION

Our approach



Signal

↓ Fourier



Mel-spectrogram

1st order derivative

↓ CNNs

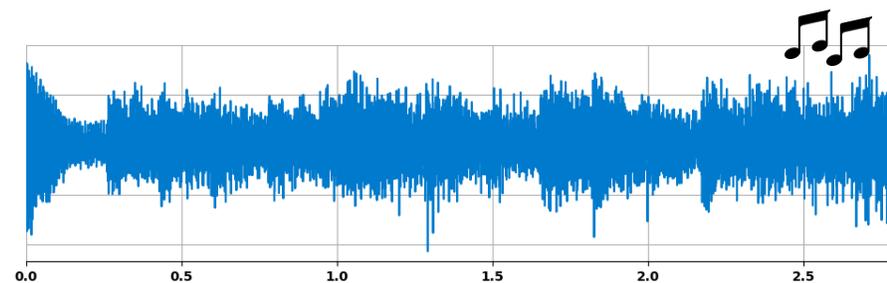
Convolutional Representations

↓ LSTM

↓ FCs

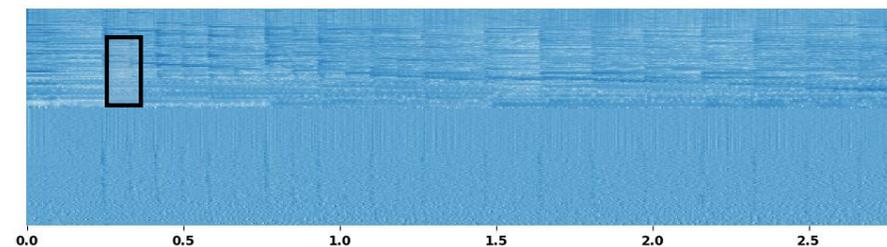
DRUM DETECTION

Our approach



Signal

↓ Fourier



Mel-spectrogram

1st order derivative

↓ CNNs

Convolutional Representations

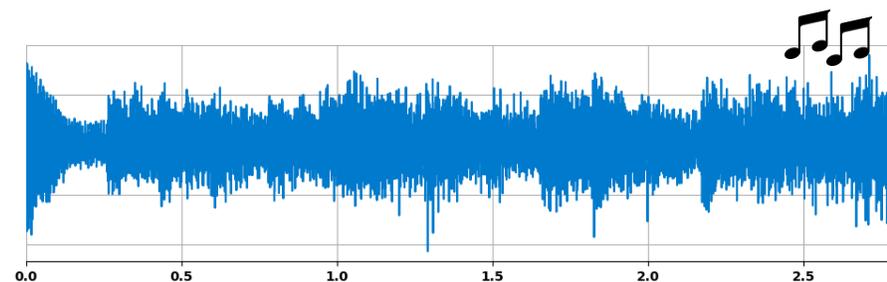
↓ LSTM

↓ FCs

Localization

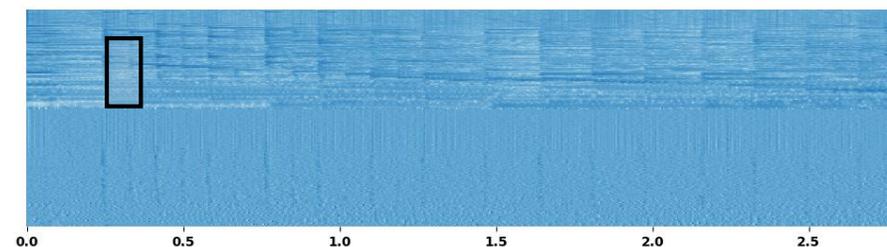
DRUM DETECTION

Our approach



Signal

↓ Fourier



Mel-spectrogram

1st order derivative

↓ CNNs

Convolutional Representations

↓ LSTM

↓ FCs

Localization

Trained with our loss
(using only occurrence counts)

DRUM DETECTION

Results

		D-DTP DATASET					
METHOD		KD	SD	HH	PRE	REC	F ₁
RANDOM	RNN	94.7	79.5	88.3	84.1	93.3	87.5
	TANHB	92.4	84.6	87.1	86.3	92.1	88.0
	RELUTS	91.3	83.8	85.2	83.7	92.3	86.8
	LSTMPB	94.4	84.1	91.4	90.8	90.8	90.0
	GRUTS	94.2	87.1	87.7	88.6	92.7	89.7
	<i>ours (LoCo)</i>	<i>92.3</i>	<i>81.2</i>	93.0	90.9	<i>87.1</i>	<i>88.9</i>
SUBSET	RNN	91.0	57.8	82.2	72.8	88.3	77.0
	TANHB	82.7	61.6	84.8	74.1	83.8	76.4
	RELUTS	79.4	62.1	80.8	69.6	84.2	74.1
	LSTMPB	85.8	68.8	83.7	78.3	84.7	79.4
	GRUTS	87.7	62.3	79.4	73.0	85.2	76.5
	<i>ours (LoCo)</i>	<i>84.9</i>	<i>59.4</i>	90.0	84.8	<i>73.5</i>	<i>78.1</i>

DRUM DETECTION

Results

		D-DTP DATASET					
METHOD		KD	SD	HH	PRE	REC	F ₁
RANDOM	RNN	94.7	79.5	88.3	84.1	93.3	87.5
	TANHB	92.4	84.6	87.1	86.3	92.1	88.0
	RELUTS	91.3	83.8	85.2	83.7	92.3	86.8
	LSTMPB	94.4	84.1	91.4	90.8	90.8	90.0
	GRUTS	94.2	87.1	87.7	88.6	92.7	89.7
	<i>ours (LoCo)</i>	92.3	81.2	93.0	90.9	87.1	88.9
SUBSET	RNN	91.0	57.8	82.2	72.8	88.3	77.0
	TANHB	82.7	61.6	84.8	74.1	83.8	76.4
	RELUTS	79.4	62.1	80.8	69.6	84.2	74.1
	LSTMPB	85.8	68.8	83.7	78.3	84.7	79.4
	GRUTS	87.7	62.3	79.4	73.0	85.2	76.5
	<i>ours (LoCo)</i>	84.9	59.4	90.0	84.8	73.5	78.1

State-of-the-art

DRUM DETECTION

Results

		D-DTP DATASET					
		KD	SD	HH	PRE	REC	F ₁
RANDOM	RNN	94.7	79.5	88.3	84.1	93.3	87.5
	TANHB	92.4	84.6	87.1	86.3	92.1	88.0
	RELUTS	91.3	83.8	85.2	83.7	92.3	86.8
	LSTMPB	94.4	84.1	91.4	90.8	90.8	90.0
	GRUTS	94.2	87.1	87.7	88.6	92.7	89.7
	<i>ours (LoCo)</i>	92.3	81.2	93.0	90.9	87.1	88.9
SUBSET	RNN	91.0	57.8	82.2	72.8	88.3	77.0
	TANHB	82.7	61.6	84.8	74.1	83.8	76.4
	RELUTS	79.4	62.1	80.8	69.6	84.2	74.1
	LSTMPB	85.8	68.8	83.7	78.3	84.7	79.4
	GRUTS	87.7	62.3	79.4	73.0	85.2	76.5
	<i>ours (LoCo)</i>	84.9	59.4	90.0	84.8	73.5	78.1

Great Overall
F1-Score

State-of-the-art



DRUM DETECTION

Results

Detection of three different drum types in drum audio extracts

Further tests on HH reveal that:



DRUM DETECTION

Results

Detection of three different **drum** 
types in drum audio extracts

Further tests on HH reveal that:

- In that setting, the **standard deviation** is only of **4.35ms** for the distance between true and predicted hits.



PIANO ONSET DETECTION

Results

Detection of **piano** notes in audio extracts 



PIANO ONSET DETECTION

Results

Detection of **piano** notes in audio extracts 

- Complex task with **88 channels**



PIANO ONSET DETECTION

Results

Detection of **piano** notes in audio extracts

- Complex task with **88 channels**
- **Tight tolerance** of 50ms for a prediction to be considered correct



PIANO ONSET DETECTION

Results

Detection of **piano** notes in audio extracts

- Complex task with **88 channels**
- **Tight tolerance** of 50ms for a prediction to be considered correct
- Comparison with **fully-supervised benchmark** models



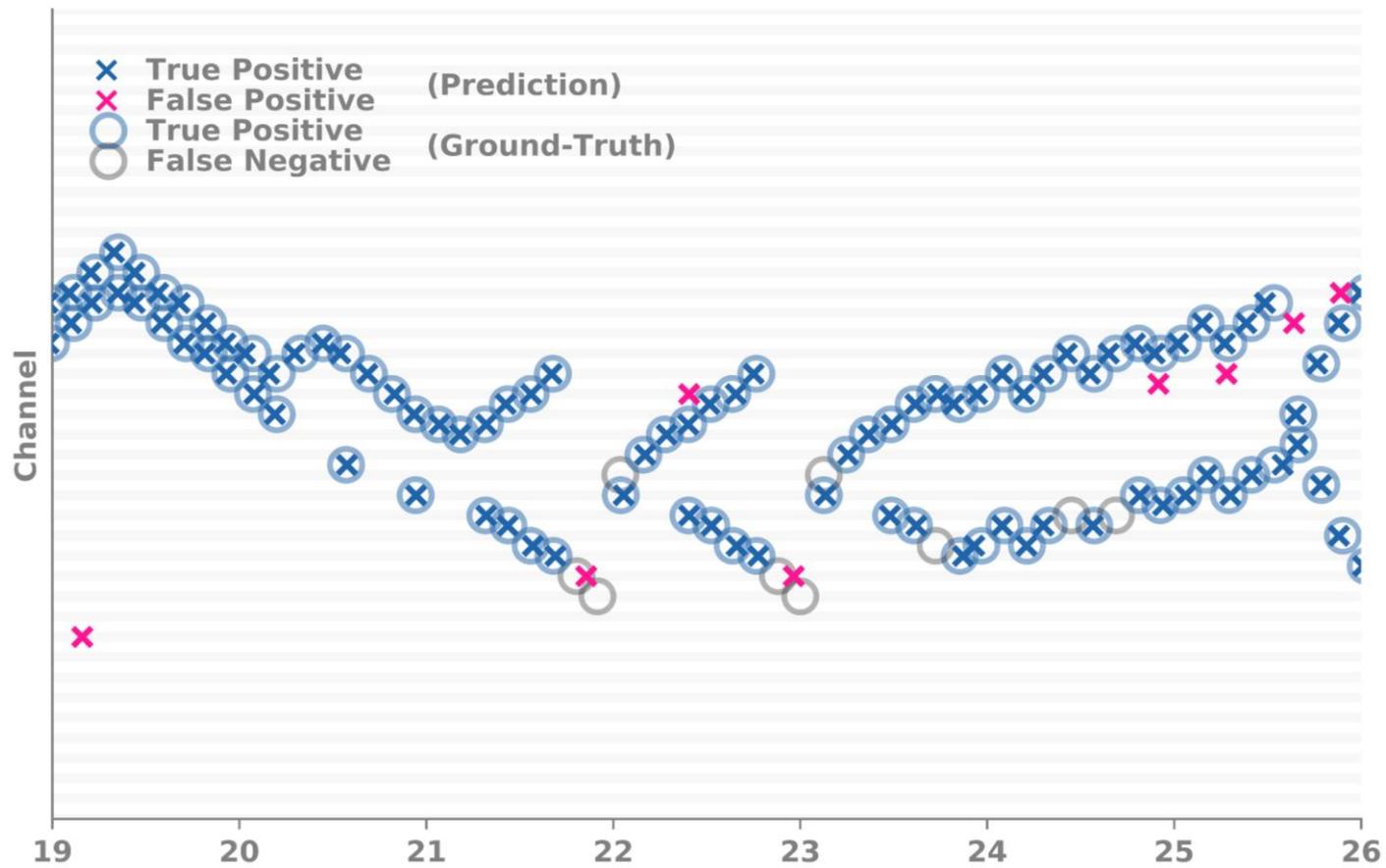
PIANO ONSET DETECTION

Results

METHOD	PRE	REC	F ₁
SIGTIA ET AL.(2016)	44.97	49.55	46.58
KELZ ET AL.(2016)	44.27	61.29	50.94
HAWTHORNE ET AL.(2017)	84.24	80.67	82.29
<i>ours (LoCo)</i>	76.22	68.61	71.99

PIANO ONSET DETECTION

Results



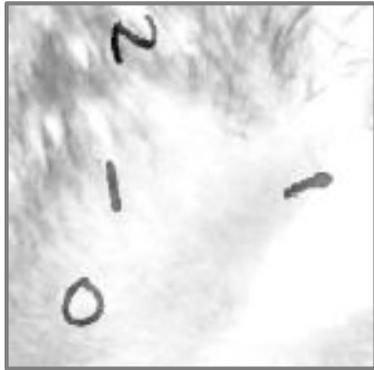
Digit Detection Experiment



DIGIT DETECTION EXPERIMENT

Main Idea

Initial Image

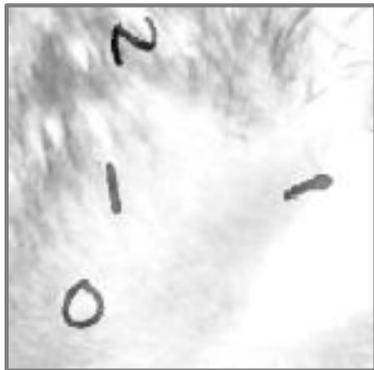


Not a sequence

DIGIT DETECTION EXPERIMENT

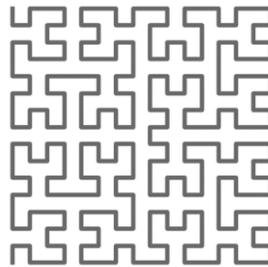
Main Idea

Initial Image



Not a sequence

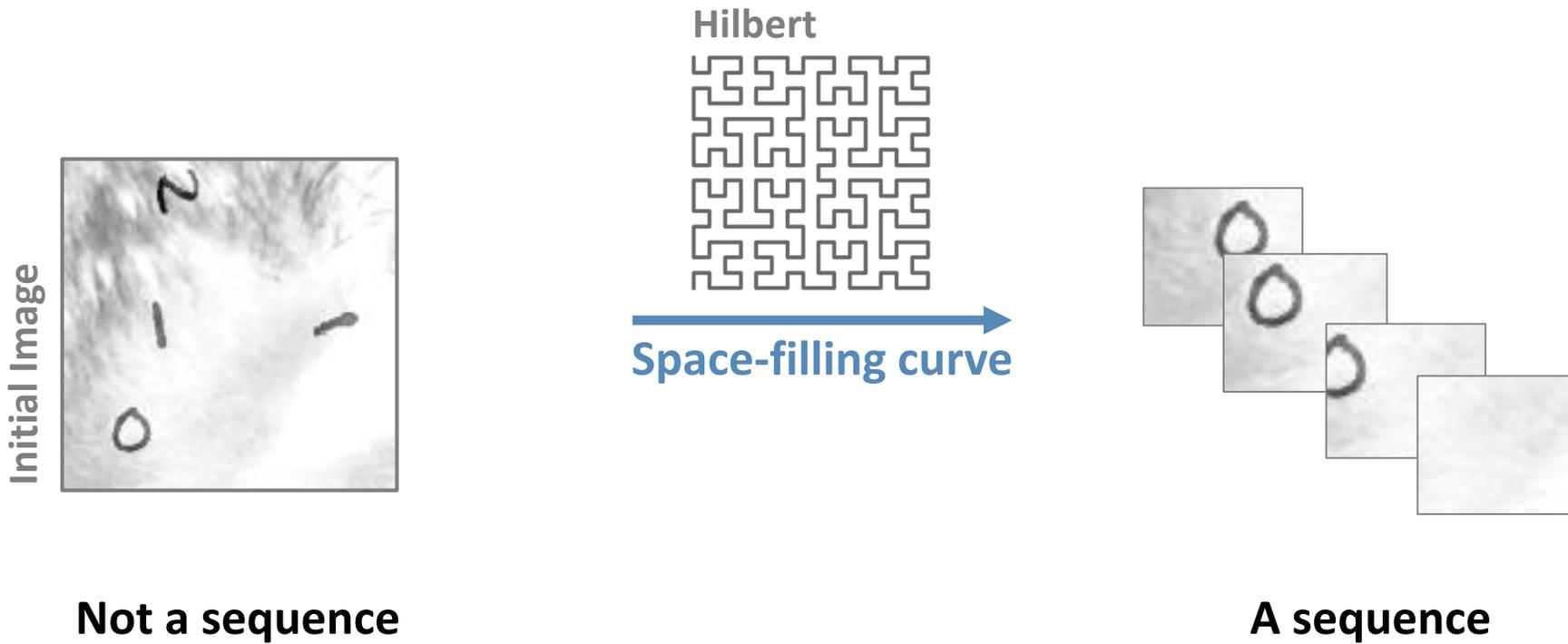
Hilbert



Space-filling curve

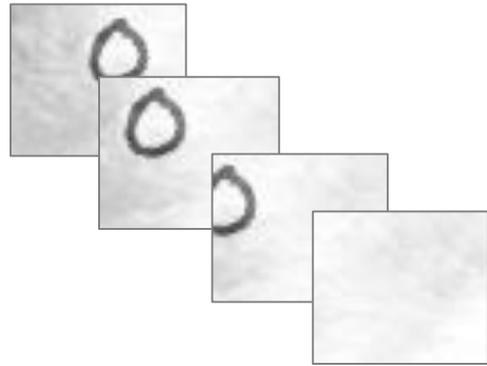
DIGIT DETECTION EXPERIMENT

Main Idea



DIGIT DETECTION EXPERIMENT

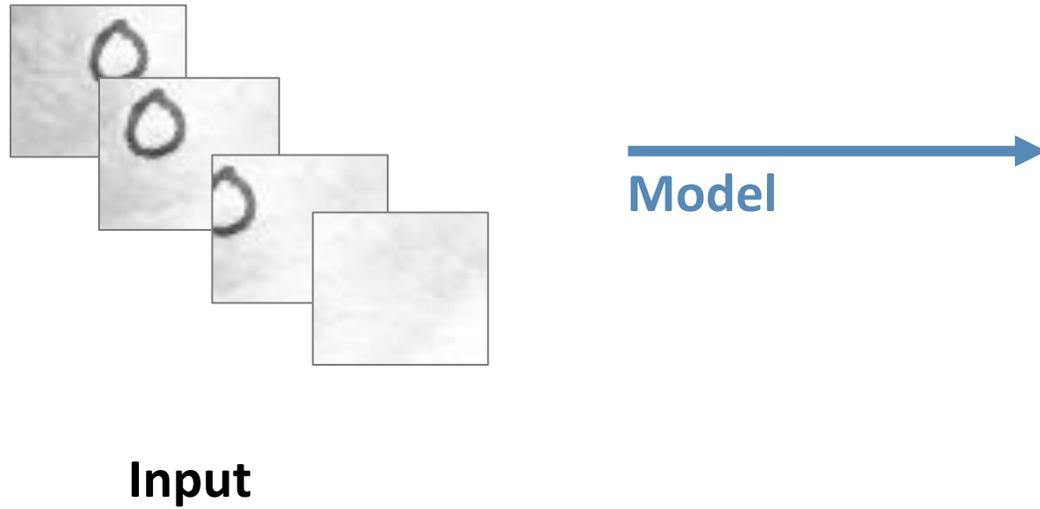
Main Idea



Input

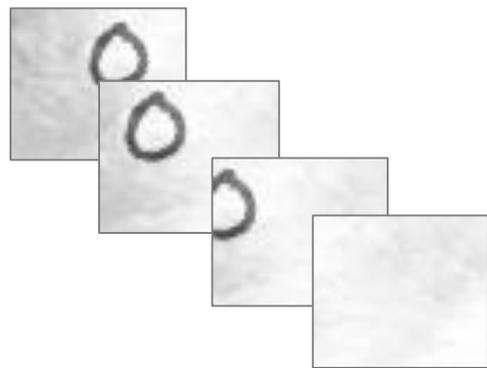
DIGIT DETECTION EXPERIMENT

Main Idea

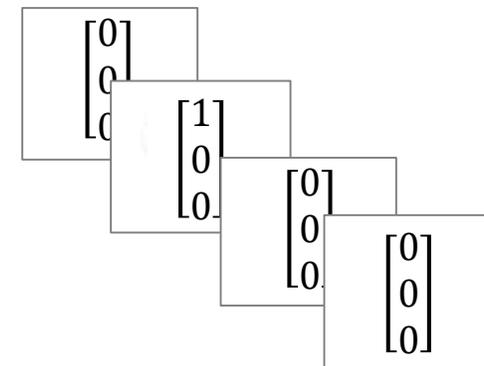


DIGIT DETECTION EXPERIMENT

Main Idea



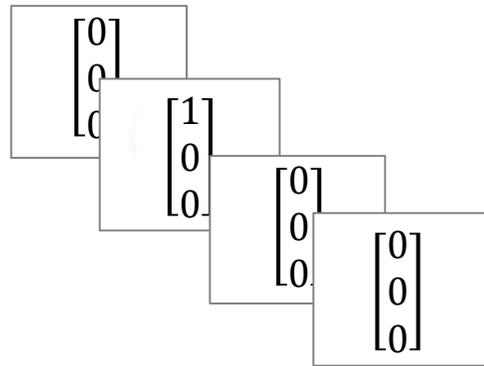
Input



Predictions

DIGIT DETECTION EXPERIMENT

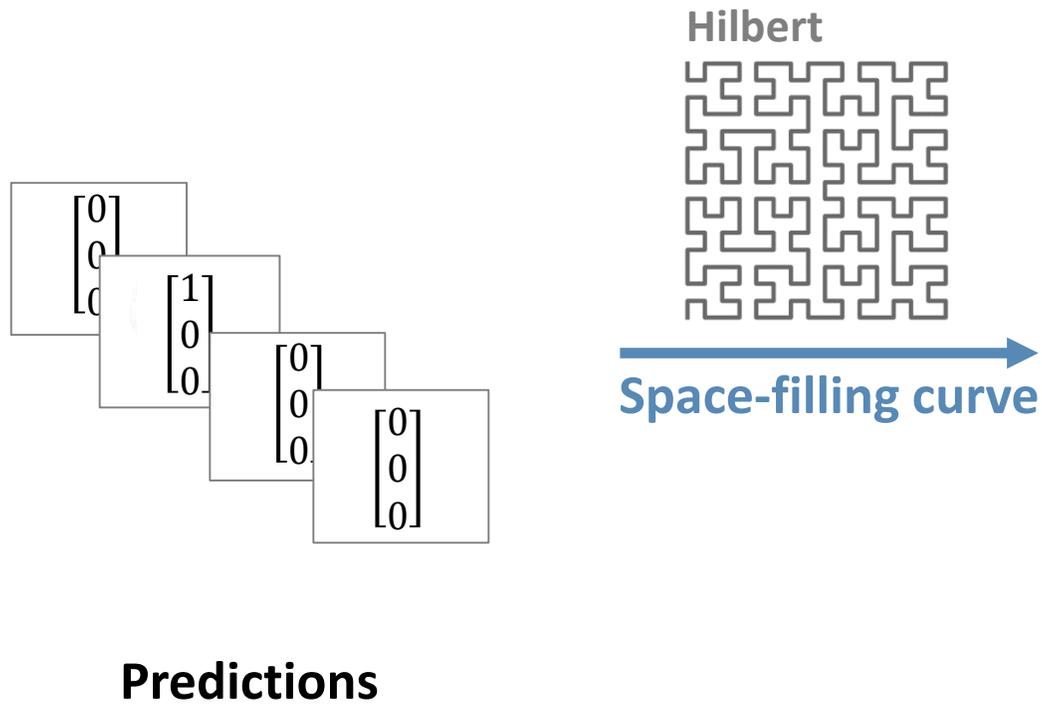
Main Idea



Predictions

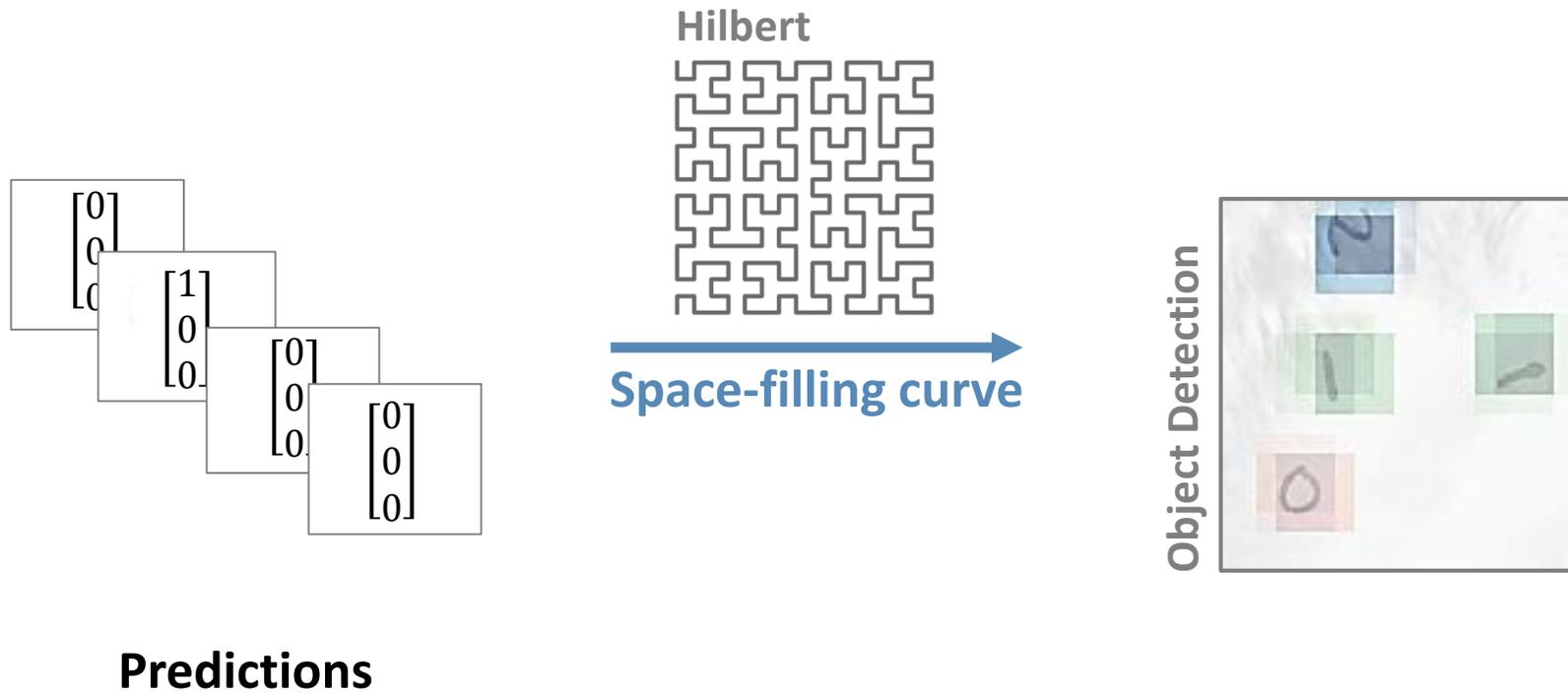
DIGIT DETECTION EXPERIMENT

Main Idea



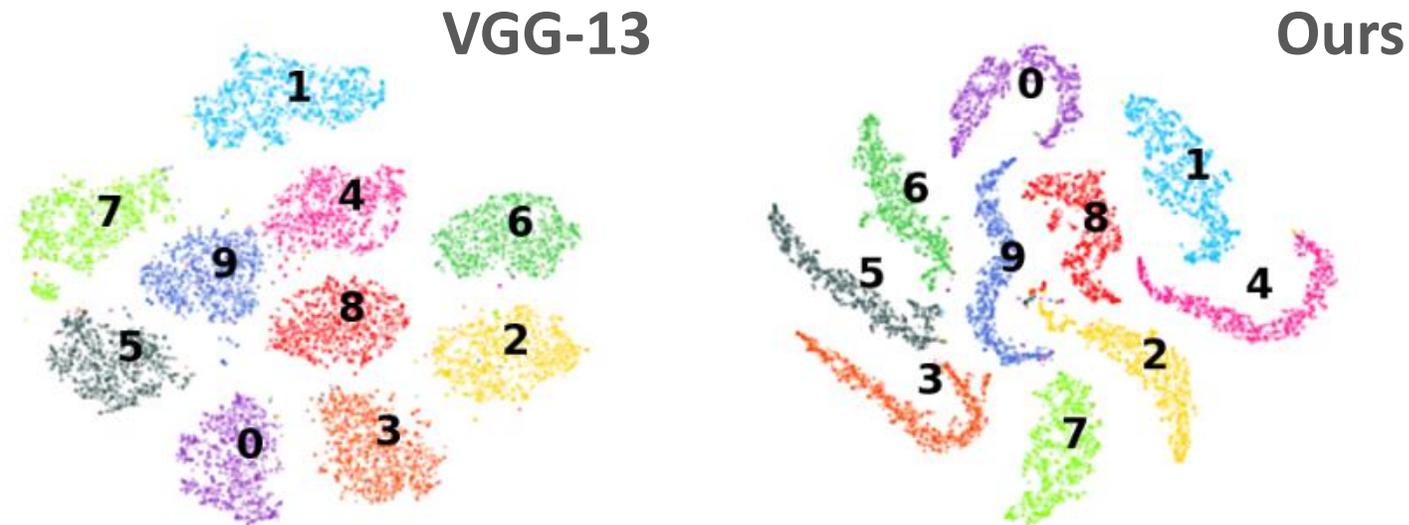
DIGIT DETECTION EXPERIMENT

Main Idea



DIGIT DETECTION EXPERIMENT

Representations

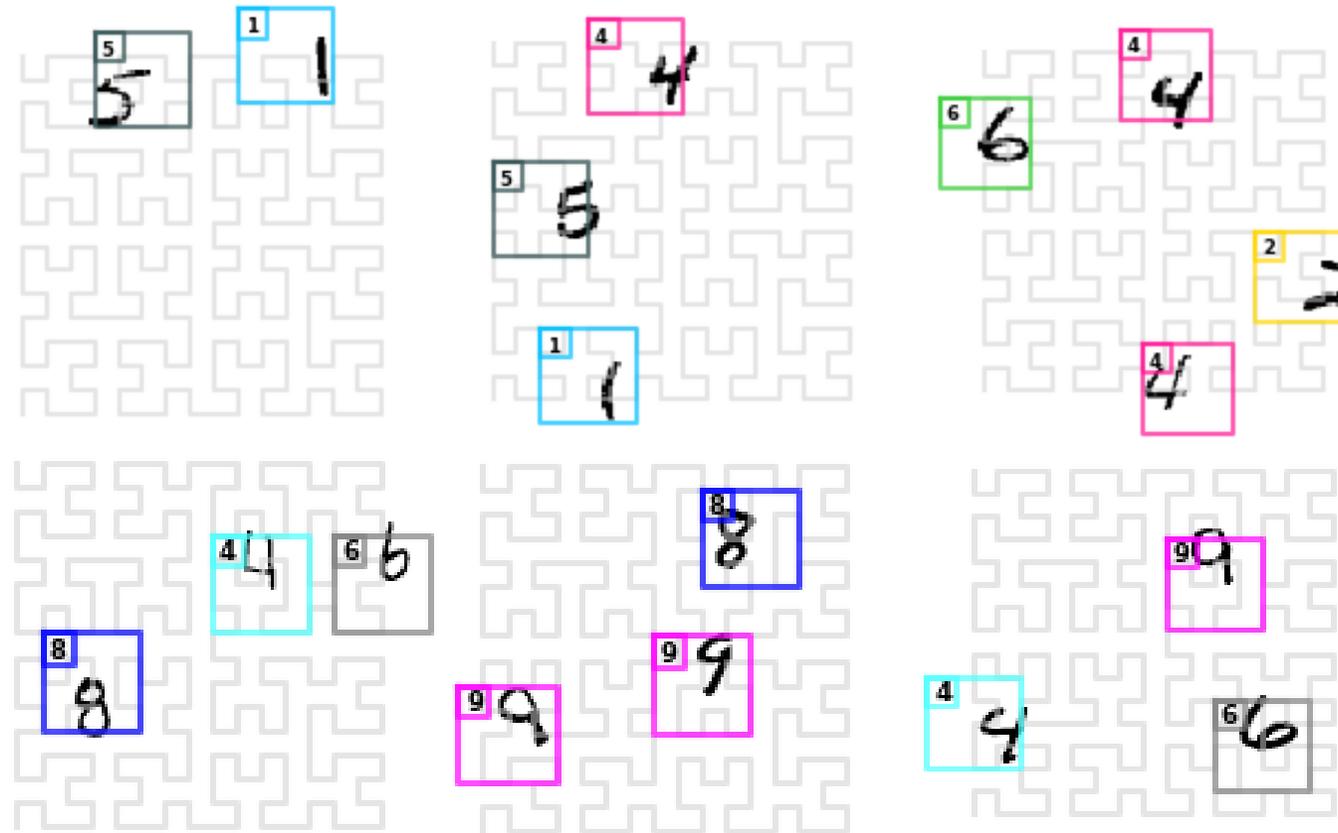


*Figure 4. Digit Representations. Comparison of t-SNE digit feature representations resulting from the *fully*-supervised VGG-13 architecture (left) and from our *weakly*-supervised approach (right).*

DIGIT DETECTION EXPERIMENT

Detection Performance

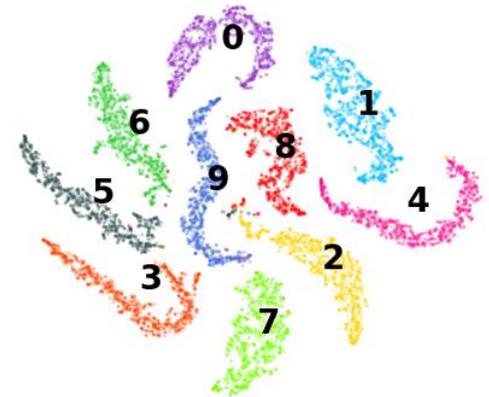
Mean absolute distance between true and estimated bounding box centers: 9:04 pixels (approx. step size of the space filling curve)



DIGIT DETECTION EXPERIMENT

Conclusion

The model learnt:

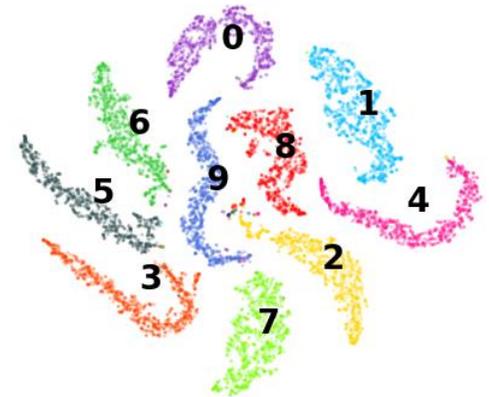


DIGIT DETECTION EXPERIMENT

Conclusion

The model learnt:

1. Feature representation

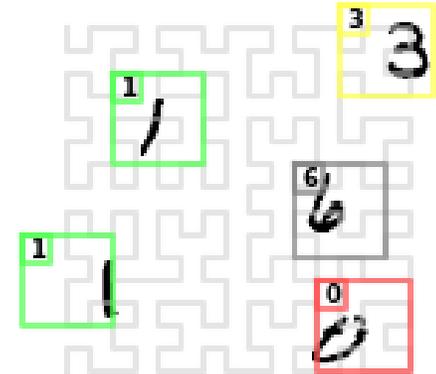


DIGIT DETECTION EXPERIMENT

Conclusion

The model learnt:

1. Feature representation
2. Space-mapping

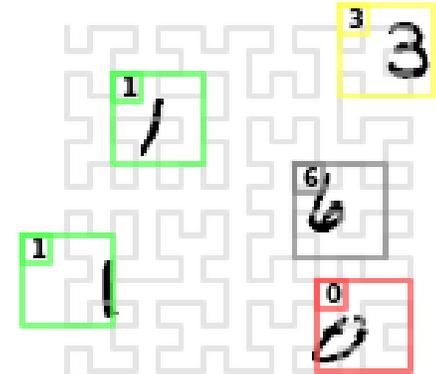


DIGIT DETECTION EXPERIMENT

Conclusion

The model learnt:

1. Feature representation
2. Space-mapping
3. Object detection

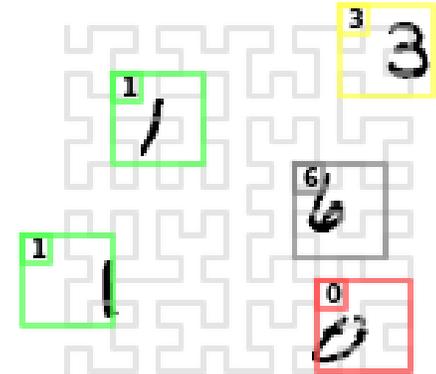


DIGIT DETECTION EXPERIMENT

Conclusion

The model learnt:

1. Feature representation
2. Space-mapping
3. Object detection



Using only occurrence counts as training labels



CONCLUSION

This work shows that implicit model constraints can be used to ensure that accurate localization emerges as a byproduct of learning to count occurrences.



CONCLUSION

This work shows that implicit model constraints can be used to ensure that **accurate localization emerges as a byproduct of learning to count occurrences.**

Competitive results against fully-supervised state-of-the-art models.

Questions?

