MusicYOLO

A Vision-Based Framework for Automatic Singing Transcription

Xianke Wang , Bowen Tian , Weiming Yang, Wei Xu and Wenging Cheng Amruta Mahendra Parulekar

• 20D070009

Keshav Singhal

• 20D070047

Sameep Chattopadhyay

• 20D070067



Some Basic Concepts



Note segmentation



Pitch extraction

5

Singing transcription



Audio event detection

Some Basic concepts

Note Segmentation

Extracting the onsets and offsets from the given audio

Segmentation necessary for subsequent pitch extraction

Two Methods- Pitch Information & Spectrogram Features

Pitch Extraction

Estimation of the Fo trajectory from the given audio signal

Traditional Methods based on time & frequency domain parameters

Current State of the Art Method: LSTM-RNN

Some Basic concepts

Singing Transcription

Challenging task even with monophonic signals

Use of audio processing algorithms similar to contour & edge detection

Current State of Art Method: RNN-HMM

Audio Event Detection

Research methods applied for both human & instruments audio detection

Traditional Methods used Support Vector Machines

Current State of the Art Method: R-CNN

Earlier AST Models

Method 1 : End-to-end approach to obtain onset, offset, and pitch probabilities. (notelevel results through post-processing)

Method 2: Onset & offset obtained first; pitch obtained through pitch extraction

Method 3: F0 tracking conducted to obtain a frame-level pitch curve; note-level onset, offset, and pitch obtained through postprocessing

MusicYOLO Method

MusicYOLO inspired by the perspective of sound event detection and based on object detection

Uses YOLOX for onset/offset detection & spectrogram peak search for pitch labeling

Detects note objects rather than isolated onset/offset moments, thus greatly enhancing the transcription performance

The MusicYOLO Model

Waveform

Pre-Processing

Note Detection

Pitch Labeling

Transcription Result



Testing Environment

Datasets



-1

ISMIR2014, MIR-ST500 Singing Transcription Dataset

Sight-Singing

Transcription

Dataset

Bach10

SSVD

Monophonic Instrument Dataset

Evaluation Metrics







Algorithms



Preprocessing

Audio Noise Reduction

Time Frequency Transformation

Linear Intensity Matching

Spectrogram Cutting

Audio Noise Reduction



Time Frequency Transformation



Linear Intensity Matching

A linear intensity mapping effectively quantifies the original spectrogram according to the spectral intensity.





Spectrogram Cutting

Spectrogram Image (Aspect ratio Conveniently process) Obtain Ratio ratio = w/h - 1where w and h represent the width and height of slice $w = frame_len \cdot scale$ $h = n_bins \cdot scale$ where $frame_len$ is the frame length of the CQT matrix,

where *frame_len* is the frame length of the CQT matrix, *n_bins* is the frequency bins of CQT, and *scale* is the scaling factor from the matrix size to image size. *frame_len* is calculated as follows:

$frame_len = t \cdot sr/hop_len$

where sr and hop_len are the CQT parameters, and t is the time length corresponding to the slice. From the above, we derive the *ratio* calculation function written as R:

 $R(t) = t \cdot \frac{sr}{hop \ len \cdot n \ bins}$



Obtain Segmentation Points

1: initialize *splits*;

2: set the iteration number r = 1 and the maximum iteration number M;

3: repeat

- 4: $last_time = splits(-1);$
- 5: $ratio = R(end_i last_time);$
- 6: **if** $ratio \ge max_ratio$ **then**
- 7: $splits = [splits, end_i];$
- 8: else if $best_ratio < ratio \le max_ratio$ then
- 9: $ratio_silence = R(start_{i+1} last_time);$
- 10: $ratio_next = R(end_{i+1} start_{i+1});$
- 11: **if** $ratio_silence > max_ratio$ **then**
- 12: $splits = C(splits, S_{end_i}, start_{i+1});$
- 13: **else if** $ratio_next > max_ratio$ then
- 14: $splits = [splits, start_{i+1}];$

15: else

16: $splits = [splits, end_i];$

17: **else**

- 18: A process similar to the above;
- 19: $r \leftarrow r + 1;$

20: **until** r > M.

Obtain approximately square slices along the spectrogram time axis



Note Detection

Note Object Detection

> Post-processing the Bounding Box

Time Shift

Note Object Detection



Note Object

Object detection to get bounding box

Left boundary = onset Right boundary = offset Bottom boundary = f_bottom

Post-processing the Bounding Box



Time Shift

- We shift the onset and offset to obtain the complete audio's onset and offset values.
- As we have previously cut the spectrogram, we splice all image slices along the time axis and add the corresponding time shift to the previously detected onset and offset.



Pitch Labelling

The Peak Search Algorithm

- The errors between F_bottom and F_pitch (actual frequency) are due to Sacrifice of frequency resolution due to time resolution and Pitch fluctuations due to vibrato and portamento
- These can be tackled using the peak search algorithm



Require: Cqt: spetrogram matrix; $onset_i, offset_i$: onset and offset of the *i*th detected note; f_{bottom_i} : bottom frequency of the *ith* detected note; *bpo*: the number of frequency bins per octave; Δb : the search range of frequency bins. **Ensure:** *Pitch*: a pitch list of the detected notes. set N as the number of the notes for each $i \in \{1, 2, ..., N\}$ do 3: $f_{init} = f_{bottom_i}$ $\Delta f 1 = f_{init} \times (1 - 2^{-\Delta b/bpo})$ 4: $\Delta f2 = f_{init} \times (2^{\Delta b/bpo} - 1)$ 5: 6: $f_{low} = f_{init} - \Delta f 1$ 7: $f_{high} = f_{init} + \Delta f2$ 8: $f_0 = 0$ 9: for each $t \in \{onset_i, \ldots, offset_i\}$ do 10: $f_0 = f_0 + f_{low} + argmax(Cqt[f_{low}: f_{high}, t])$ 11: $f_{pitch} = f_0 / (offset_i - onset_i)$ 12: $Pitch_i = 69 + 12 \times \log(f_{pitch}/440)$



Results

Results For Different Model Configuration

NOTE DETECTION RESULTS OF DIFFERENT MODEL CONFIGURATIONS

		CQT	MEL	YOLOX	Faster RCNN	SSVD P B F				ISMIR201 B	4 F	Mixed P B F		
Onset	M1 M2 M3	√ √	√	✓ ✓	√ ∕	95.03 95.25 92.33	95.65 95.28 94.79	95.32 95.24 93.45	92.40 89.85 89.53	89.39 86.46 85.20	90.71 89.85 87.00	93.23 94.37 91.44	92.00 90.45 89.86	92.50 92.21 90.43
Offset	M4 M1 M2 M3 M4	√ √		✓ ✓	✓ ✓ ✓	92.06 97.49 97.72 94.85 94.64	94.57 98.13 97.75 97.41 97.25	93.19 97.78 97.71 96.02 95.81	89.87 87.10 88.87 85.53 85.72	84.60 84.53 81.88 81.65 80.93	86.85 85.65 85.04 83.26 82.97	91.15 92.34 93.66 90.97 90.81	91.32 99.96 89.58 89.25	90.03 91.73 91.62 90.07 89.79



Why YOLOX ?

- The decoupling detection heads make the model converge faster and perform better.
- The use of anchor free for training eliminates the process of clustering to obtain prior bounding boxes from the dataset, which reduces the overfitting potential.
- SimOTA enables the model to automatically analyze the number of positive samples corresponding to each ground truth.

Note Detection Results

				NOTE DI	ETECTION	RESULTS (OF DIFFERI	ent Mode	ELS				
		ISMIR2014			MST500			SSVD			Bach10		
		Р	R	\mathbf{F}	Р	R	\mathbf{F}	Р	\mathbf{R}	\mathbf{F}	Р	\mathbf{R}	\mathbf{F}
	Tony	72.13	64.19	67.63	50.98	56.88	53.45	79.00	78.96	78.90	91.94	88.55	90.15
Onset	Fu&Su	83.00	75.43	78.76	50.78	53.97	52.10	79.70	81.04	80.11	87.17	79.41	83.00
	Wang I	81.09	82.49	81.60	73.50	78.59	75.67	76.88	82.44	79.31	87.38	88.65	87.61
	MusicYOLO I	89.55	82.76	85.88	81.89	75.02	78.17	88.02	86.88	87.41	92.56	84.68	88.31
	Wang II	88.38	40.15	53.90	65.07	29.70	40.34	91.36	65.79	75.70	84.74	24.48	37.21
	MusicYOLO II	91.01	89.24	90.01	62.53	65.20	63.59	95.99	96.38	96.17	93.97	94.19	94.07
	Tony	79.88	70.37	74.47	55.82	62.44	58.59	93.59	93.50	93.43	95.71	92.16	93.84
	Fu&Su	79.60	72.94	75.87	53.26	56.82	54.74	84.48	86.09	84.99	92.39	84.07	87.91
Offset	Wang I	76.39	77.86	76.95	65.25	69.89	67.25	71.86	77.78	74.43	87.00	88.70	87.43
	MusicYOLO I	83.43	77.07	79.99	78.78	72.01	75.10	85.14	84.21	84.64	97.24	88.80	92.68
	Wang II	21.39	9.12	12.46	26.50	12.05	16.38	13.84	10.00	11.45	42.47	12.33	18.80
	MusicYOLO II	85.79	84.34	84.96	62.89	65.57	63.94	97.89	98.28	98.07	96.80	97.06	96.91

Note Detection Error Analysis

• We classify the detection errors into four categories-



Note Detection Error Analysis





Transcription Results

		1	ISMIR201	1 KANSC	MST500			SSVD			Bach10		
		Р	R	F	Р	R	F	Р	R	F	Р	R	F
	Tony	65.69	58.60	61.70	38.92	43.29	40.78	71.54	71.44	71.41	91.07	87.74	89.31
	Fu&Su	77.81	70.62	73.80	36.46	38.55	37.33	63.08	64.07	63.37	83.93	76.48	79.92
ote	Wang I	66.79	68.15	67.33	66.09	70.88	68.18	58.24	62.25	60.00	82.12	83.65	82.55
ž	MusicYOLO I	85.24	78.84	81.78	74.69	68.61	71.42	79.19	78.15	78.64	89.27	81.72	85.20
	MusicYOLO II	84.00	82.44	83.12	45.16	47.08	45.94	85.77	86.10	85.92	91.77	91.99	91.86
	Tony	49.99	44.87	47.10	24.98	27.98	26.27	67.12	67.16	67.07	84.65	81.77	83.13
fse	Fu&Su	62.52	56.87	59.40	22.69	24.03	23.25	56.96	57.60	57.17	77.79	70.96	74.12
e w/ of	Wang I	53.09	54.09	53.49	46.61	50.04	48.12	44.36	47.24	45.64	73.77	75.04	74.19
	MusicYOLO I	72.21	66.90	69.34	61.22	56.27	58.56	68.62	67.77	68.17	87.28	79.91	83.30
Vot	MusicYOLO II	72.18	70.96	71.49	31.11	32.54	31.70	84.45	84.78	84.60	89.28	89.49	89.37

Transcription Error Analysis

• We classify the transcription errors into two categories-



TOTAL

ERRORS

Transcription Error Analysis



- The figure shows that the onset detection errors are the main contributor to transcription error.
- Meanwhile pitch errors are rarely found.
- This indicates that extracting pitch in singing transcription is easy
- Onset/offset detection is the key to transcription performance.



THANK YOU

This Photo by Unknown author is licensed under <u>CC BY</u>.