
Automatic Music Transcription

Ojas Chaturvedi, Kayshav Bhardwaj, Tanay Gondil, Ritwik Jayaraman, Franklin Shang, Emily Li, Sean Su, Shreeya Sarurkar, Elliott Soderberg, Arnav Kalekar, Vishaal Iyer, Junyong Lee, Basil Khwaja

- AI4Musicians Lab



Input:
music audio signal

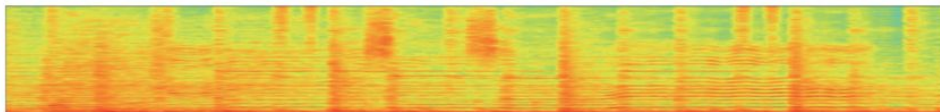


Input:
music audio signal

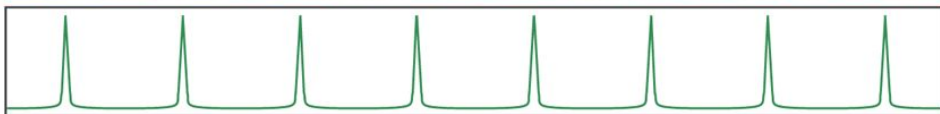


Short-time Fourier transform & beat tracking

Spectrogram



Tatum times

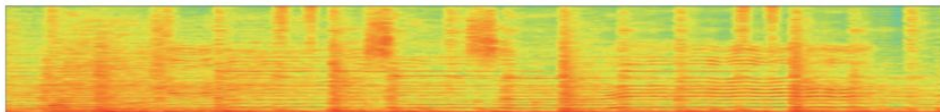


Input:
music audio signal

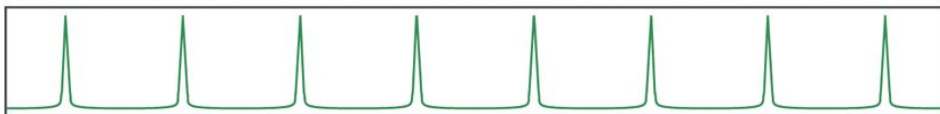


▼ Short-time Fourier transform & beat tracking

Spectrogram



Tatum times



▼ Proposed method

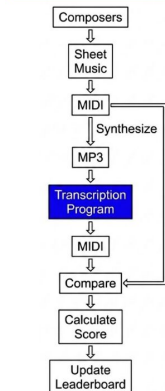
Output:
musical score



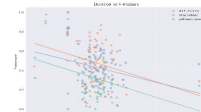
Past and Current Work

Advancing Multi-Instrument Music Transcription: Results from the 2025 AMT Challenge

Ojas Chaturvedi (ochaturv@purdue.edu), Kayshav Bhardwaj, Tanay Gondil, Benjamin Shiue-Hal Chou, Yujia Yan, Kristen Yeon-Ji Yun, Yung-Hsiang Lu, Sungkyun Chang



Performance Considerations

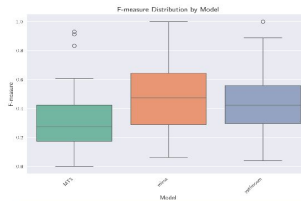


We observed two clear trends: model performance declined on longer pieces and on pieces with more overlapping instruments.

We launched the 2025 AMT Challenge with a novel multi-instrument dataset, revealing a critical gap in polyphonic transcription.

Results

Rank	Model Name	F1 Score	Precision	Recall	Overlap	Runtime
1	MIROS	0.5998	0.6558	0.5724	0.7391	22.05
2	YourMT3-YPTF-MoE-M	0.5938	0.6010	0.5888	0.7305	12.60
3	YourMT3-YPTF-S	0.5581	0.5565	0.5615	0.7326	15.40
4	YourMT3-P	0.3947	0.3966	0.3985	0.7263	14.99
5	MT3 (baseline)	0.3932	0.3911	0.4115	0.7180	20.19
6	YourMT3-YPTF-SP-V	0.3305	0.3280	0.3358	0.7147	14.50
7	press_to_win 1	0.1199	0.1105	0.1346	0.7331	19.30
8	press_to_win 2	0.1190	0.3094	0.3331	0.7310	18.08
9	YourMT3-YPTF-MoE-MP	0.2173	0.2150	0.2206	0.6116	16.03
10	press_to_win 3	0.2168	0.2144	0.2203	0.6159	16.15
11	ByteDance Piano	0.1721	0.2041	0.1689	0.5423	9.67
12	press_to_win 4	0.1470	0.1305	0.1799	0.6998	21.74
13	ReconVAT	0.1415	0.1215	0.1803	0.7908	5.45
14	Basic Pitch	0.0634	0.0550	0.0782	0.5977	3.91



Architectures

- MIROS**
 - Uses MusicFM, a conformer-based, self-supervised foundation model, which is pretrained via BEST-RQ masked token modeling to leverage abundant unlabeled audio data.
 - Extends the YourMT3 framework with parallel 15-style multi-decoders, utilizing rotary position embeddings and hardware-optimized FlashAttention.
 - Showed better accuracy on the new competition data compared to systems potentially overfit to existing datasets.
- YourMT3-YPTF-MoE-M**
 - Achieved competitive accuracy (F1 0.5938) with a significantly fast runtime (12.60 ms), demonstrating excellent practical efficiency.
 - Achieved the highest Recall (0.5888) among all submissions, indicating superior success in detecting all true musical notes.
 - Demonstrated robust foundational performance on single-instrument tracks (F1 0.7584).
 - Part of a model family that incorporates enhancements like Mixture-of-Experts (MoE) routing within its sequence-to-sequence architecture.



Dataset

The evaluation set consists of 76 newly composed pieces – including modern atonal works and underrepresented instruments such as bassoon and viola – rendered from MIDI using FluidSynth (Piaf3 GM). Each piece follows strict constraints: tempo between 60-90 BPM, meters limited to 3/4, 4/4, or 6/8, and rhythmic subdivisions no smaller than sixteenth notes. To avoid ambiguity, elements like swing, trills, double-dotted rhythms, and grace notes were excluded. Pieces span a pitch range of C2-C7, use dynamics from *pp* to *ff*, and include up to three instruments drawn from eight allowed choices, with at most one string instrument per piece.

Scoring

Two MIDI files are compared to calculate the transcription program's score. The evaluation metrics include precision, recall, F1 score, and overlap ratio. The precision and recall of the reference and estimated MIDI are computed using the *mt_eval* library. Overlap is calculated through onset and offset, i.e., the timing of a music note's beginning and ending, by computing the intersection over union between transcription and ground truth.



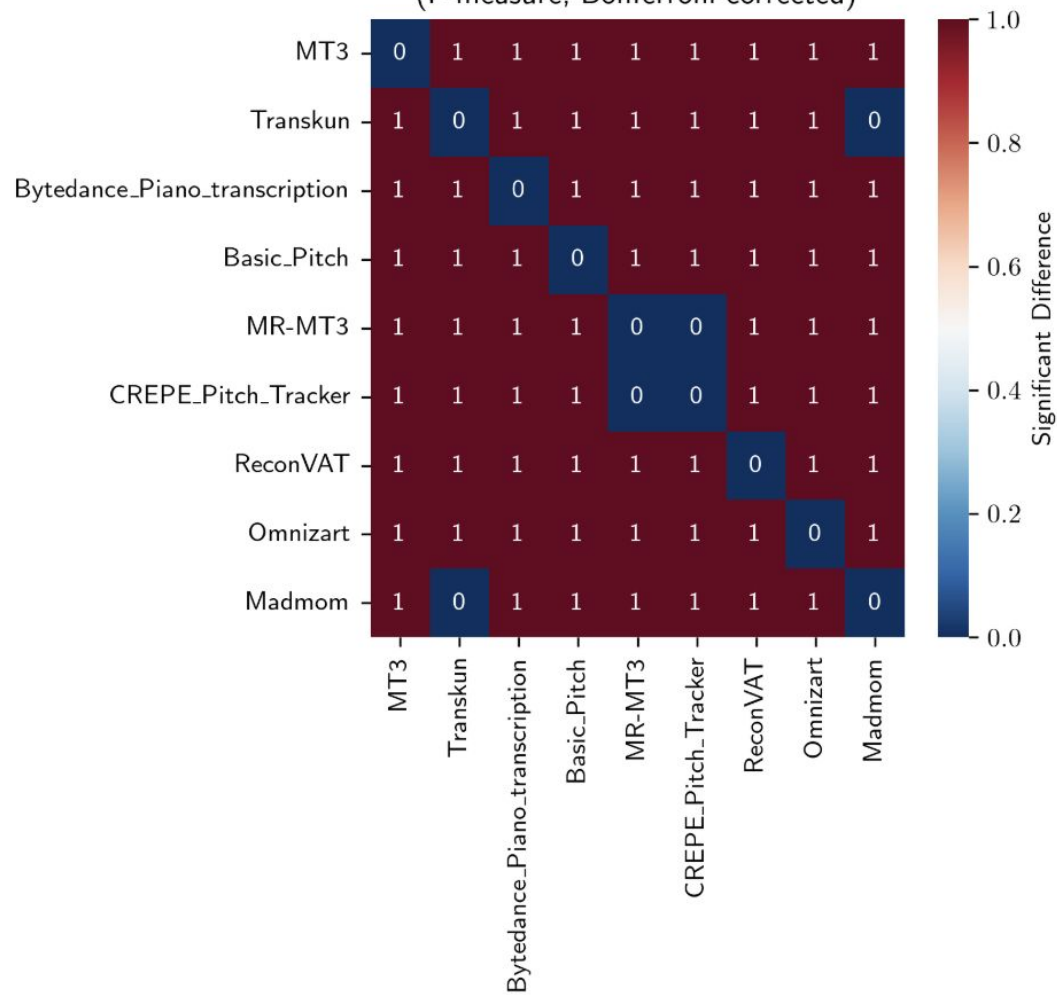
Future Work

Our goal is to build an **interpretable** AMT model – one that approaches state-of-the-art accuracy while offering greater transparency. We aim to design a system whose internal decisions can be examined and understood, rather than treated as a black box.

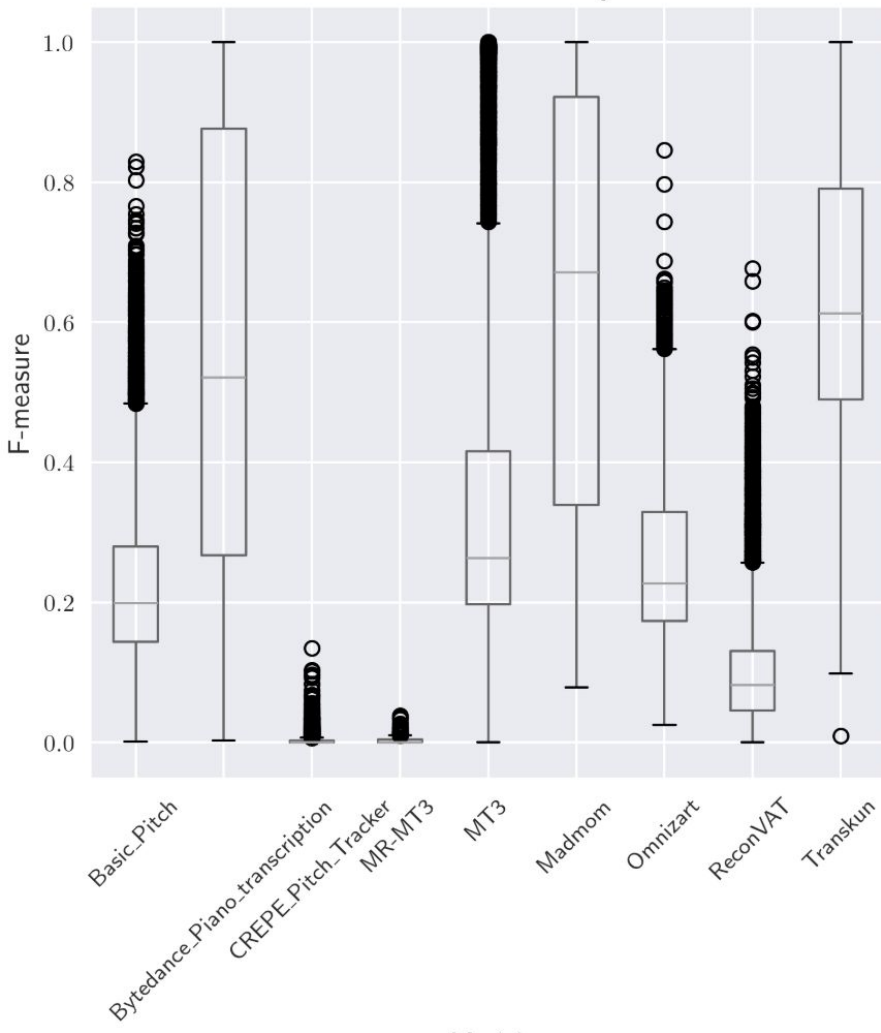
To achieve this, we plan to employ the DDS (Differentiable Dictionary Search) architecture. DDS uses a normalizing flow to learn a dictionary of possible "notes" and represents an audio signal as a linear combination of these learned bases. This allows the model not only to make accurate predictions but also to reveal which dictionary elements it relies on, providing a clearer window into how the transcription is produced.



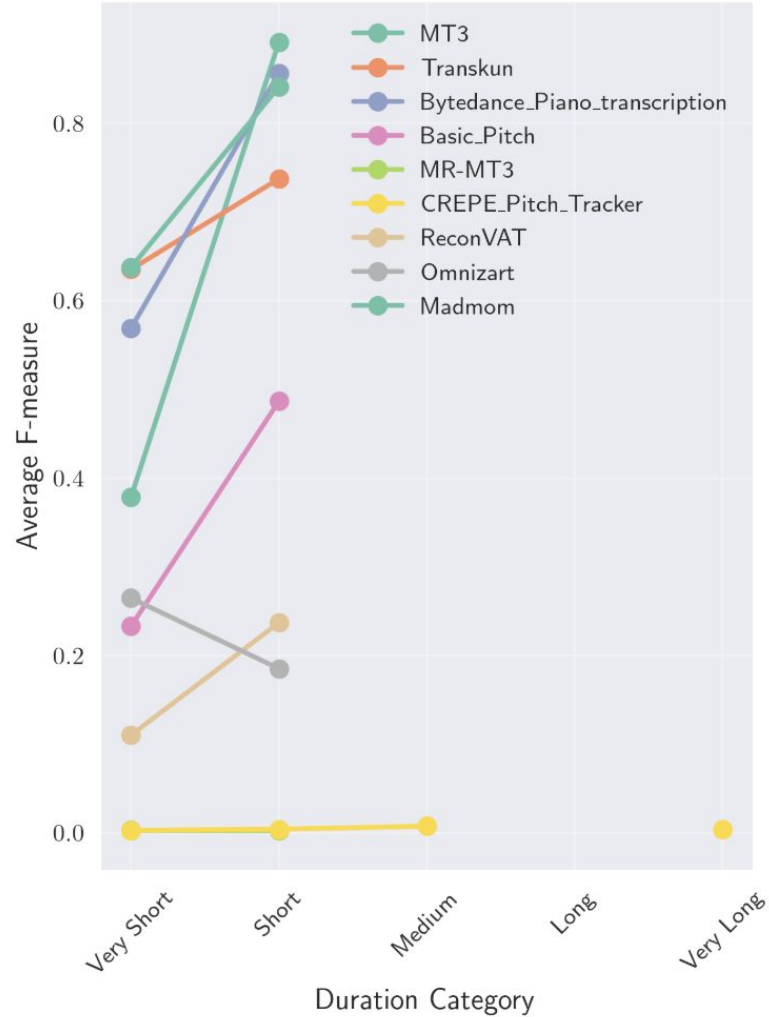
Significant Model Differences
(F-measure, Bonferroni-corrected)



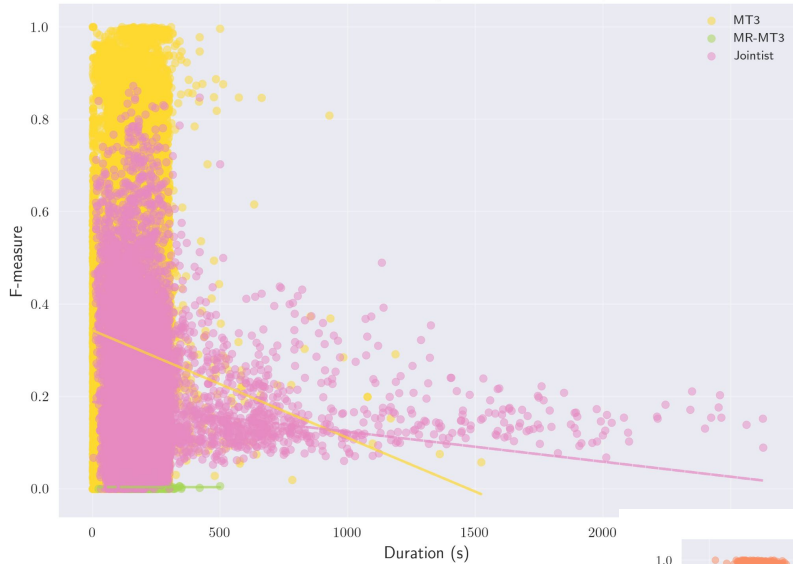
F-measure Distribution by Model



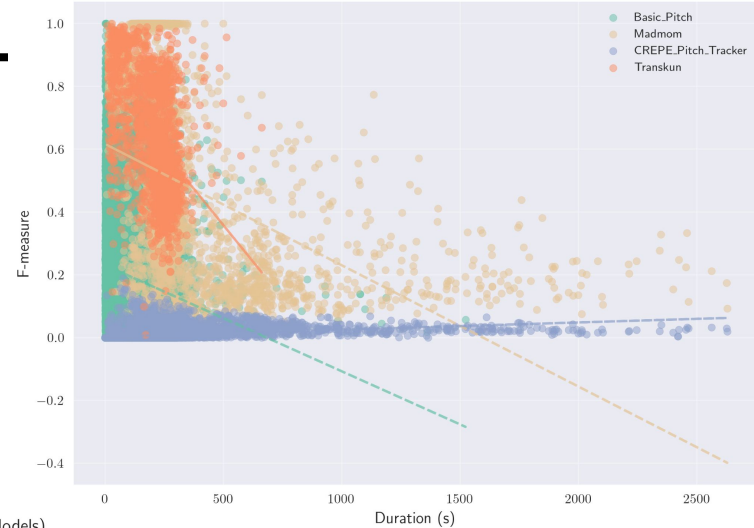
Performance vs Duration Category by Model



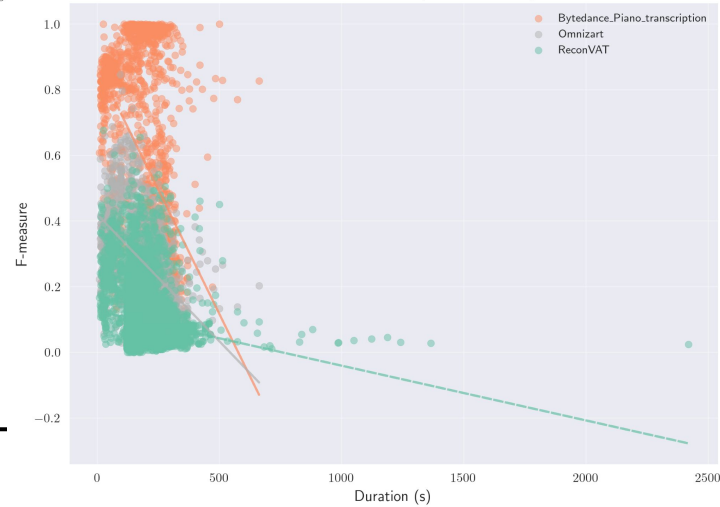
Duration vs F-measure (Transformer Models)



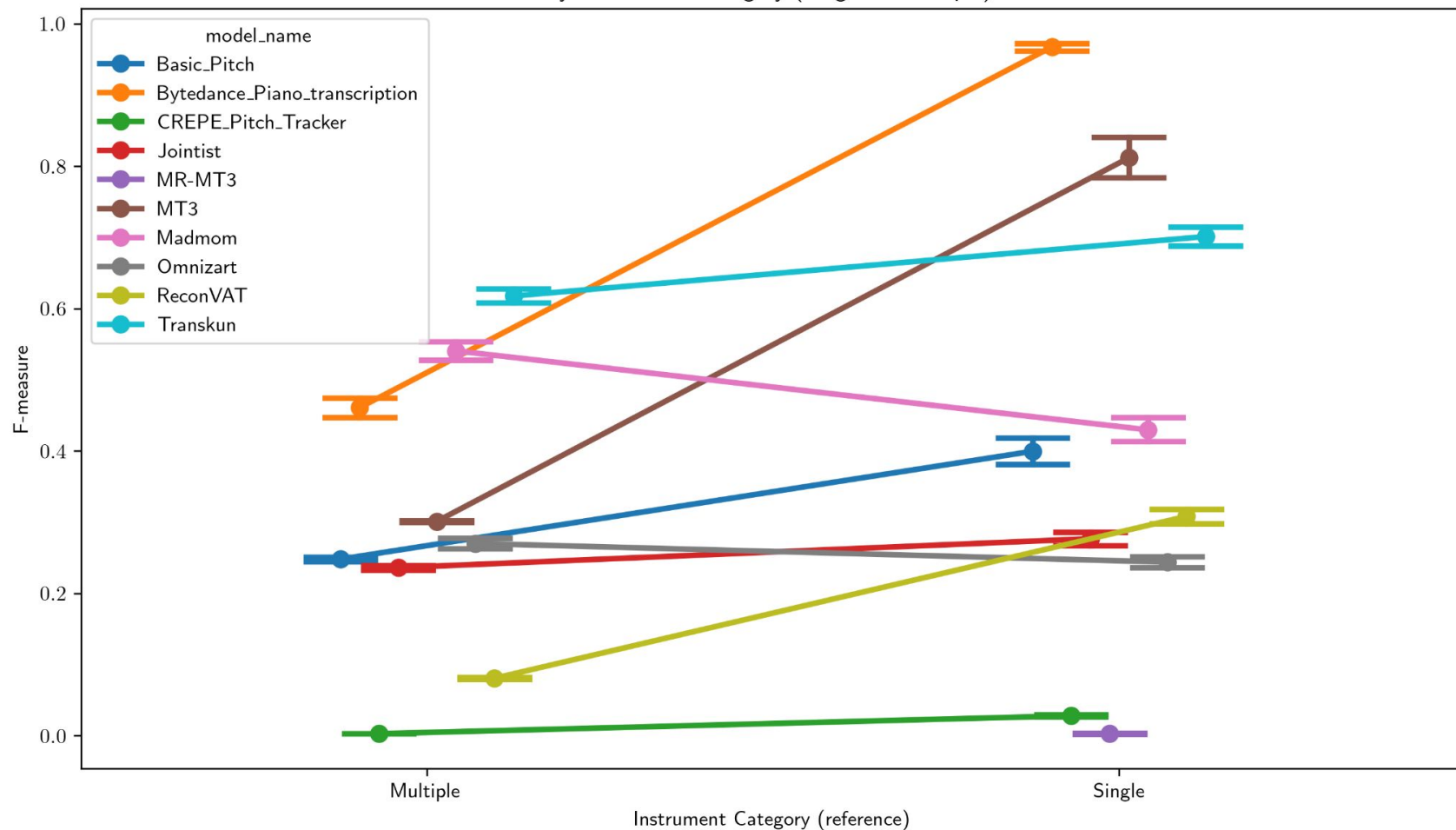
Duration vs F-measure (Traditional Models)



Duration vs F-measure (Neural Models)



F-measure by Instrument Category (Single vs Multiple) and Model



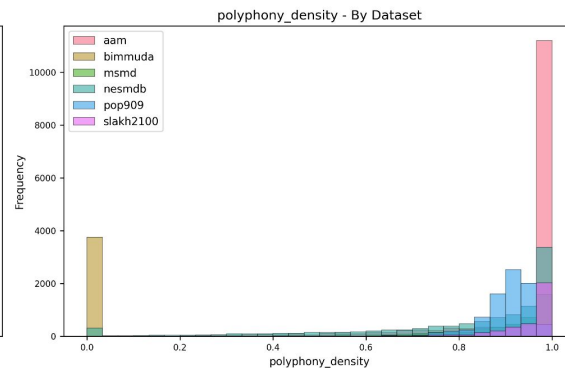
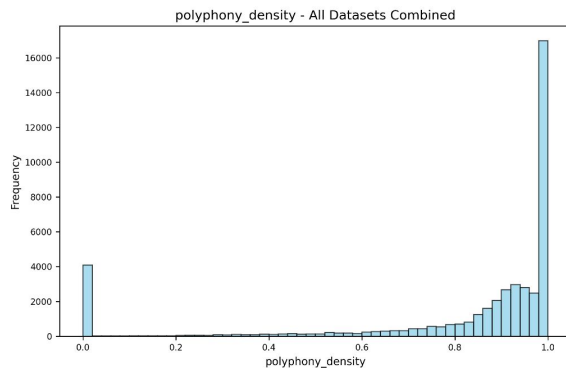
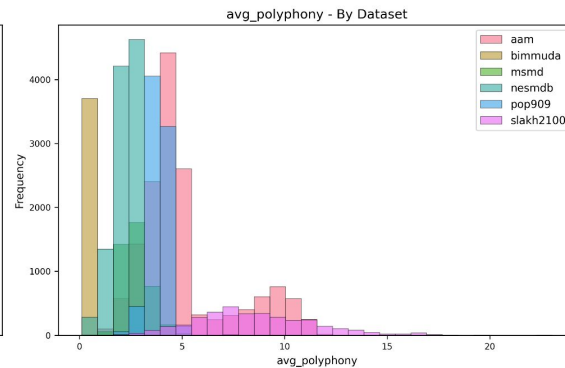
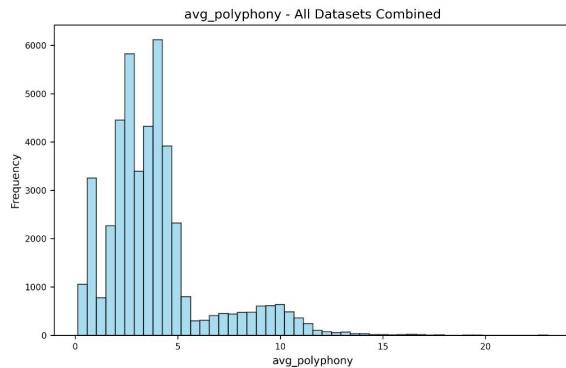
Music Complexity - Metrics

The image displays three musical examples illustrating different complexity metrics:

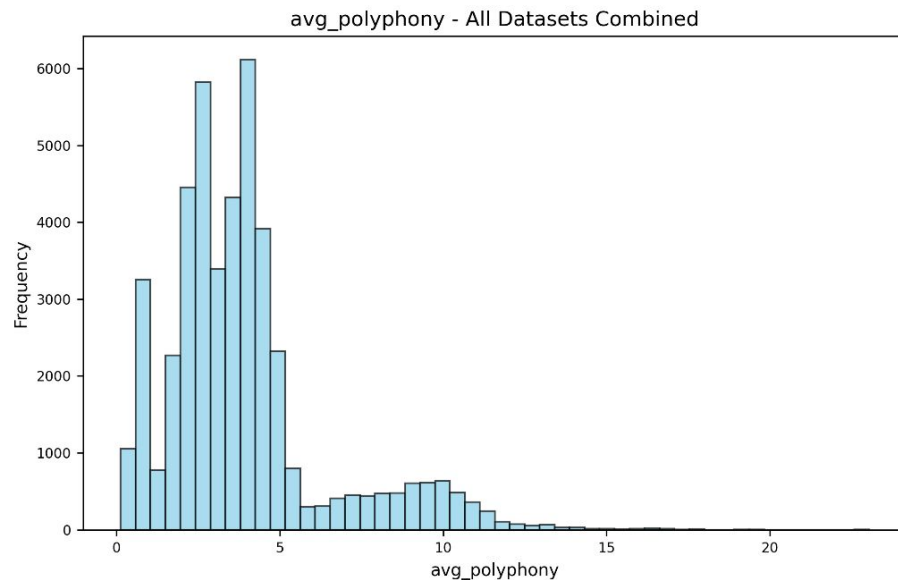
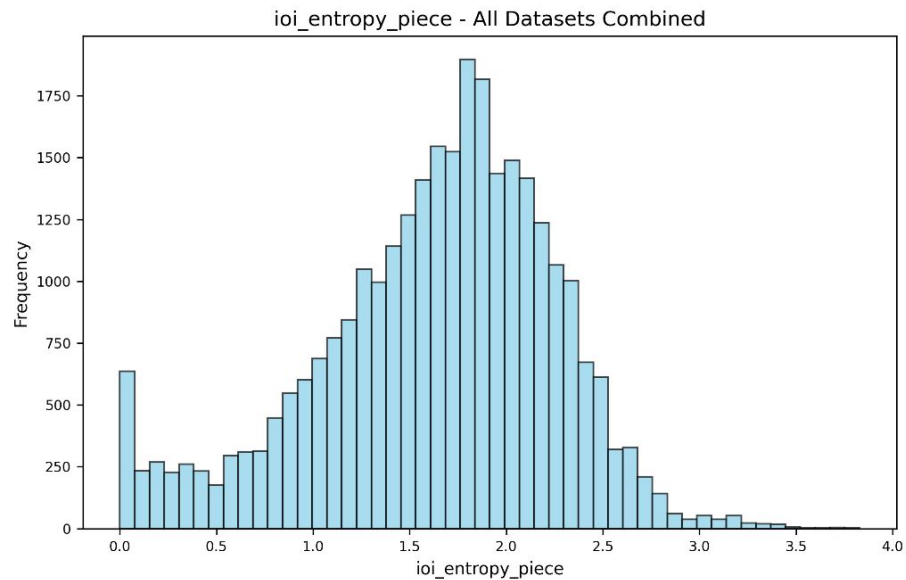
- Top Example (Rhythmic):** Shows a melody and harmony in 4/4 time. The melody has a blue circle around the first two notes. The harmony consists of chords: Dm (9th), F (7th), Dm (7th), G (6th), and Am (R). A red oval highlights the first two chords (Dm and F).
- Bottom Example (Polyphonic):** Shows a melody and harmony in 4/4 time. The melody has a blue circle around the first two notes. The harmony consists of chords: F (9th), Dm (11th), Em (9th), Dm (11th), and Em (R). A red oval highlights the first two chords (F and Dm).

- Rhythmic
 - Hypothesis: Higher variability in note placements (less common rhythms) → worse transcription accuracy
 - Example: Entropy of Inter-Onset-Intervals
- Harmonic
 - Hypothesis: More complex chord transitions → worse transcription accuracy
 - Example: Average Transition Complexity
- Polyphonic
 - Hypothesis: Higher number of notes played simultaneously → worse transcription accuracy
 - Example: Proportion of piece with 2+ notes being played simultaneously

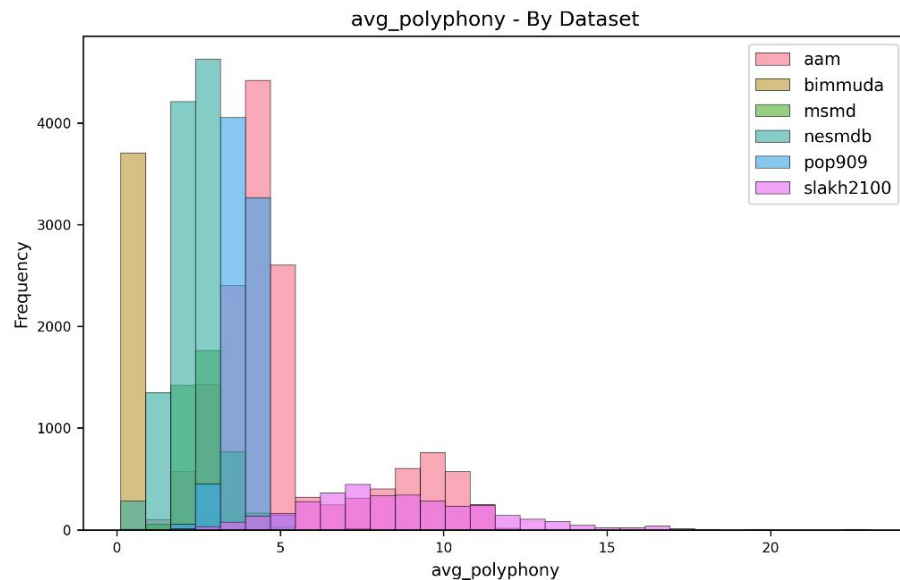
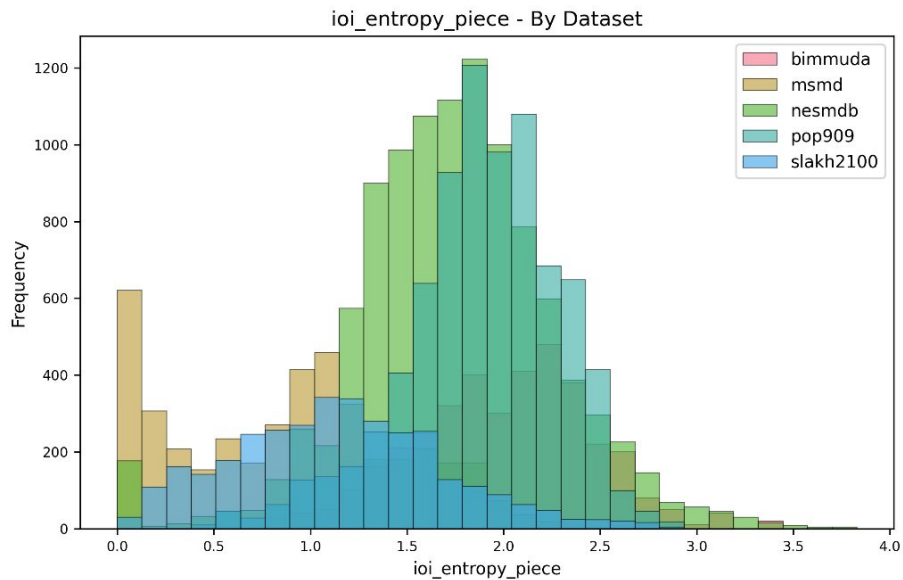
Complexity Results



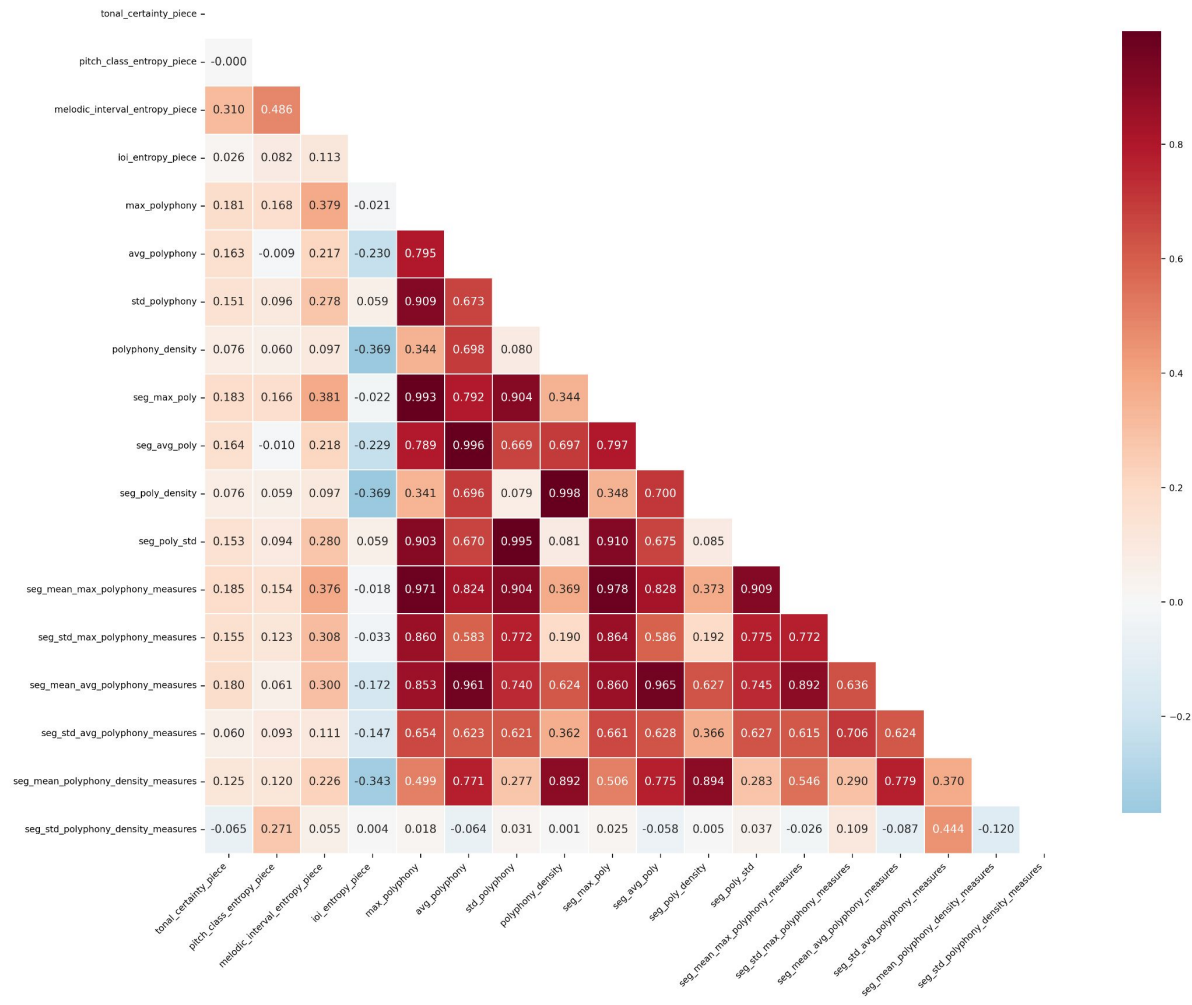
– Complexity Results



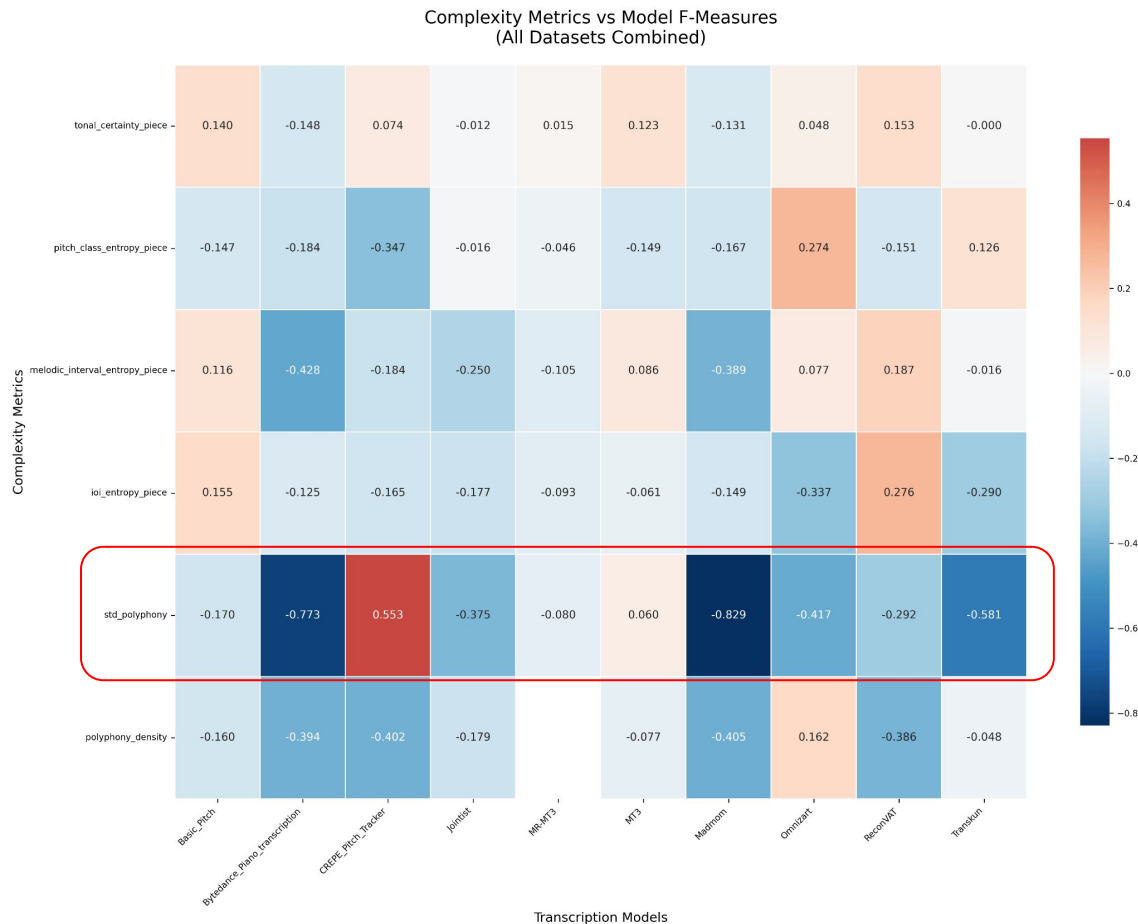
– Complexity Results



Complexity Metrics Intercorrelations
(All Datasets Combined)



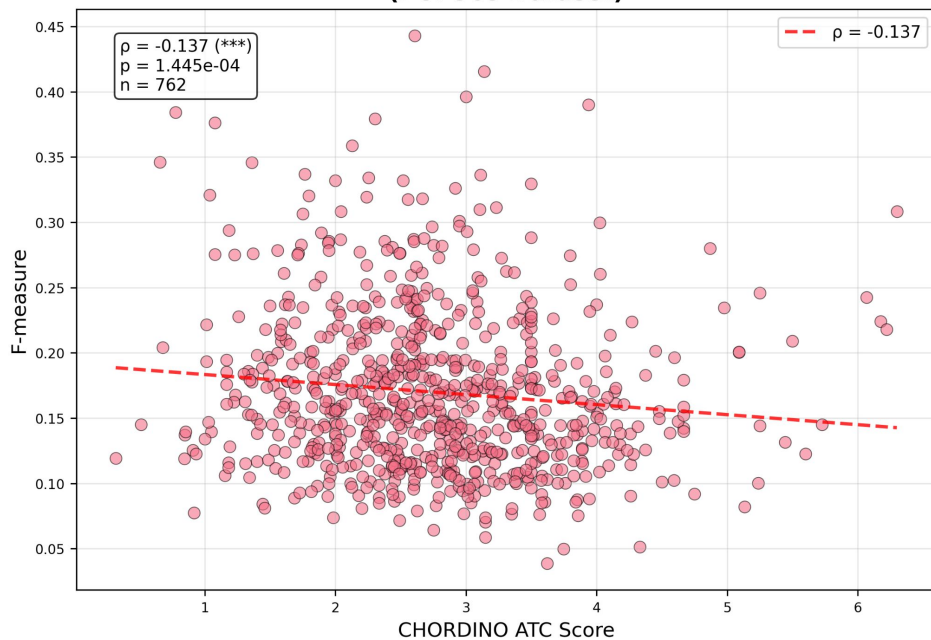
– Complexity Results



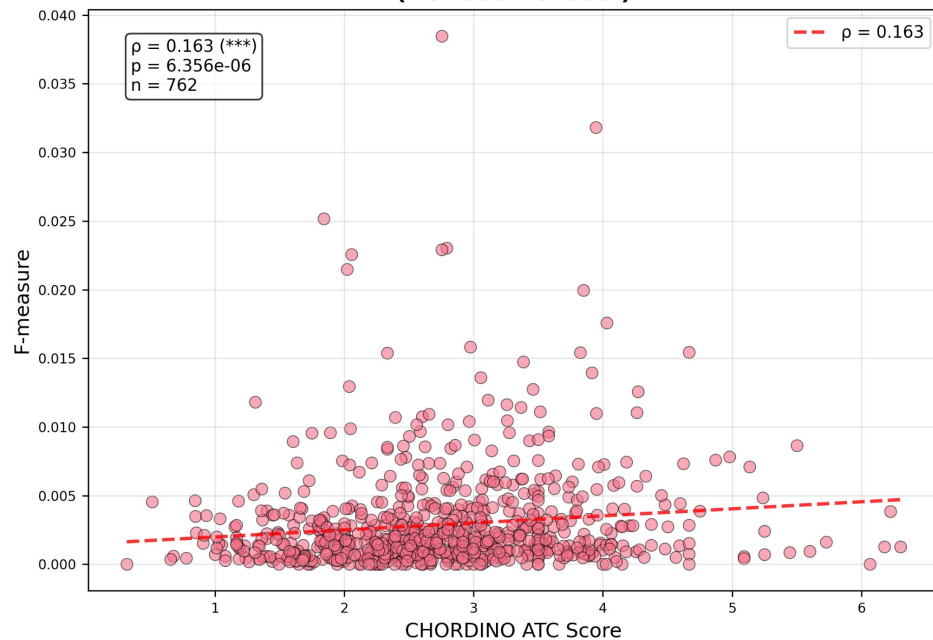
- Positive correlation
 - Higher complexity score -> higher accuracy
- Negative correlation
 - Higher complexity score -> lower accuracy
- Weak correlations
 - Could correspond to the amount of data models were trained on

– Complexity Results

**Jointist: CHORDINO ATC vs F-measure
(POP909 Dataset)**

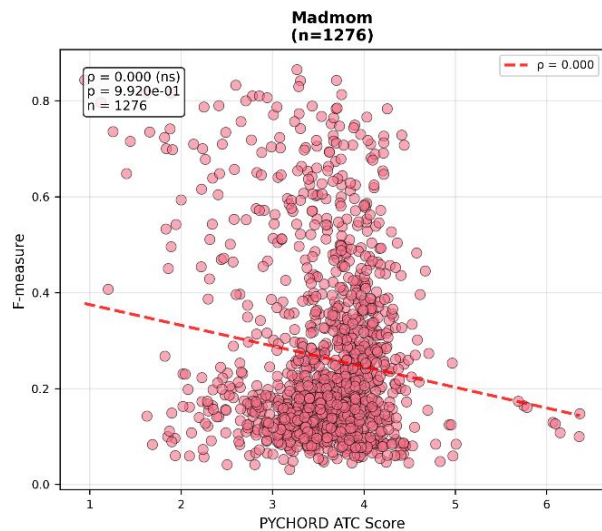
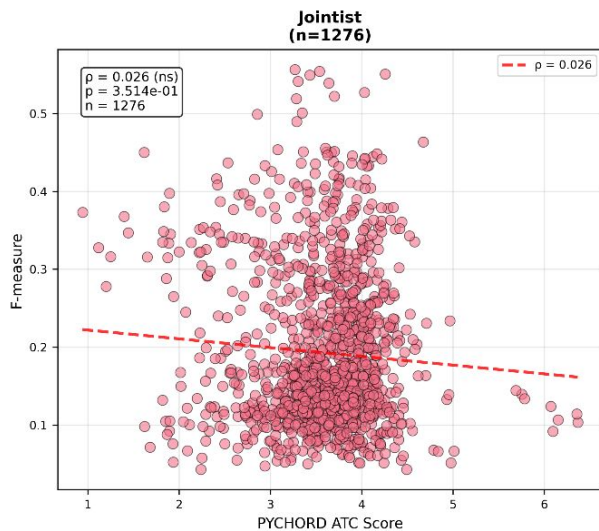
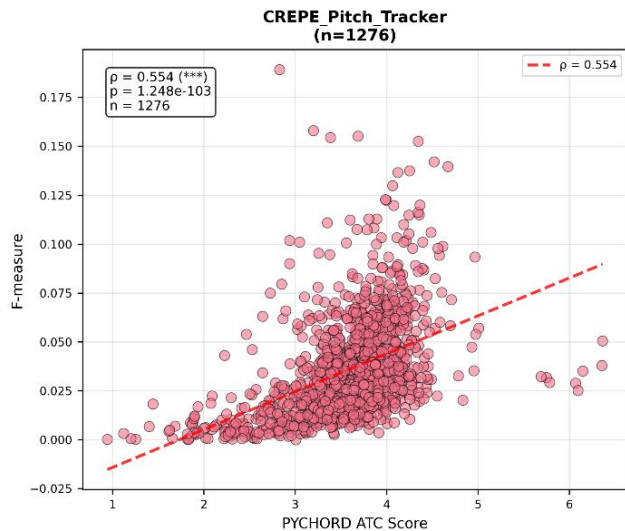


**CREPE_Pitch_Tracker: CHORDINO ATC vs F-measure
(POP909 Dataset)**

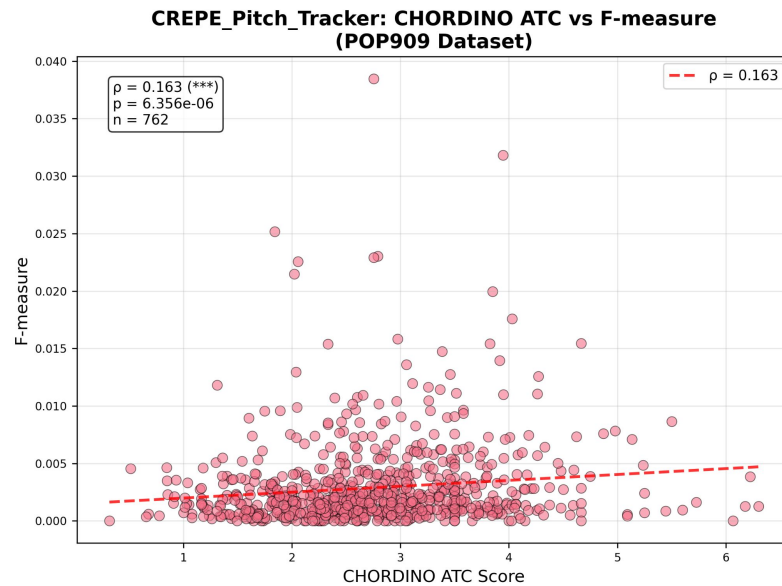
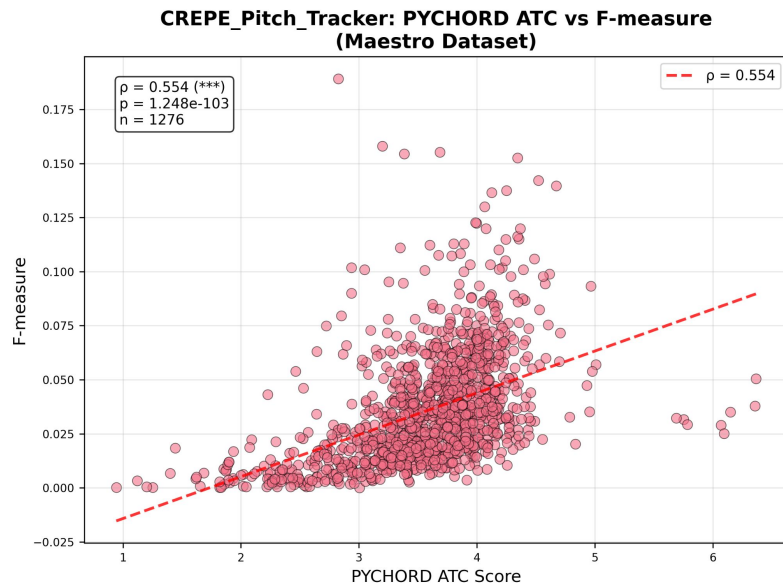


— Complexity Results

**PYCHORD ATC vs F-measure by Model
(Maestro Dataset)**



— Complexity Results

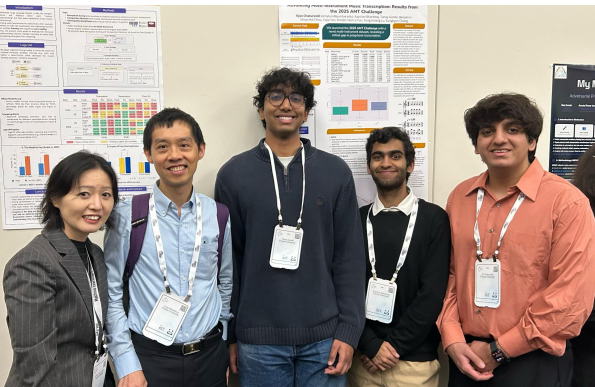


Results Analysis

- Positive correlation
 - Higher complexity score -> higher accuracy
 - Negative correlation
 - Lower complexity score -> higher accuracy
 - Weak correlations
 - Could correspond to the amount of data models were trained on
-

Future Work

Research Dissemination



The 2025 Automatic Music Transcription Challenge

- Competitors from 5 countries
- Presented results at the AI for Music Workshop at NeurIPS

Paper in progress conducting a literature review in SotA AMT, aiming to publish in IEEE MultiMedia



Drafting a paper to be published to Arxiv

Building an Interpretable Model

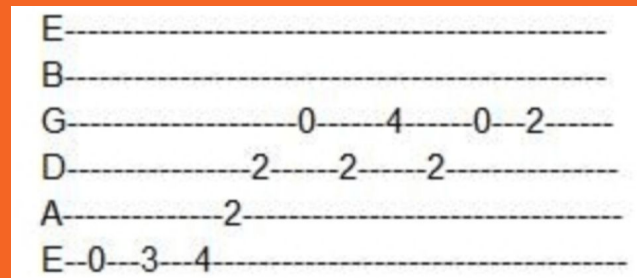
- What is interpretability?
 - Current models are “black boxes”
 - We don’t know what they “learn”
 - Errors are difficult to isolate and fix
 - The overfitting problem



Building an Interpretable Model

- Our goal is to build an interpretable AMT model:
 - Approaching SoTA accuracy with higher transparency
 - Look inside the model to see why it does what it does
 - Planned DDS architecture
 - Differentiable Dictionary Search
 - Deep learning to build a dictionary of possible “notes”
 - Sums and generates the predicted linear sum of the bases
-

Computer Vision



Converting guitar performance video to tablature

TAB staff lines represent guitar strings

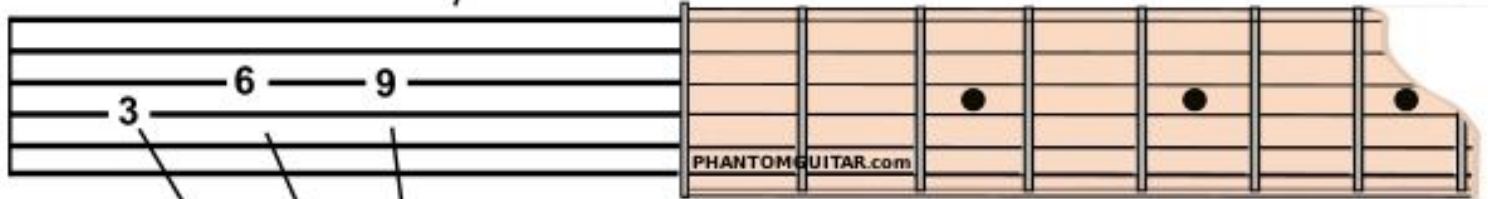
String numbers

1
2
3
4
5
6

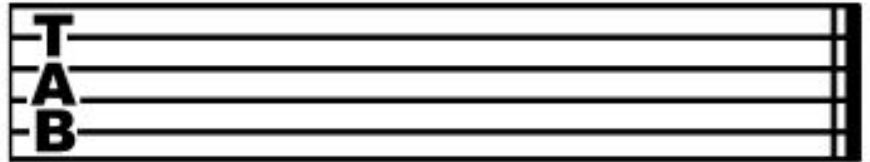
e
B
G
D
A
E

String names

Fret numbers placed on the staff
tell you where to put your fingers.

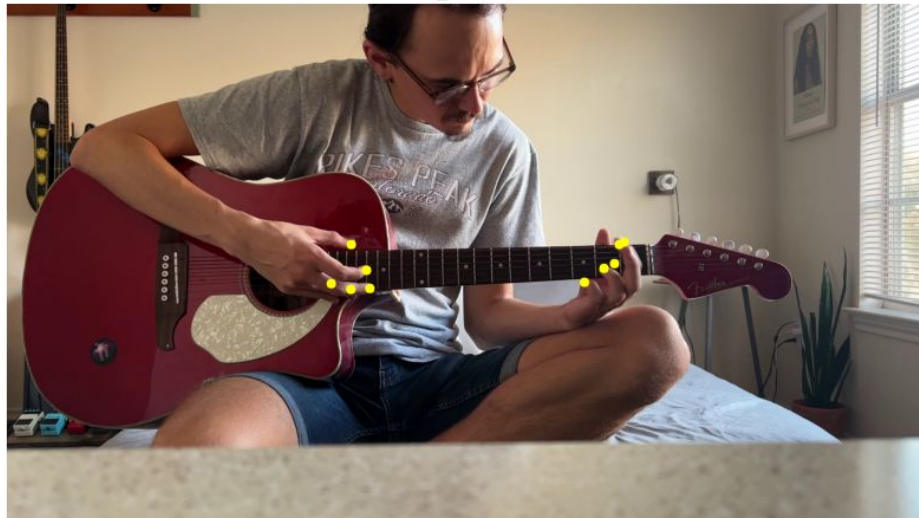


The word **TAB** is placed at the start
of a TAB staff so it is not confused
with regular music notation..



Fingertip Detection (MediaPipe)

31_14.png



34_8.png



Fretboard Isolation (Mask R-CNN)

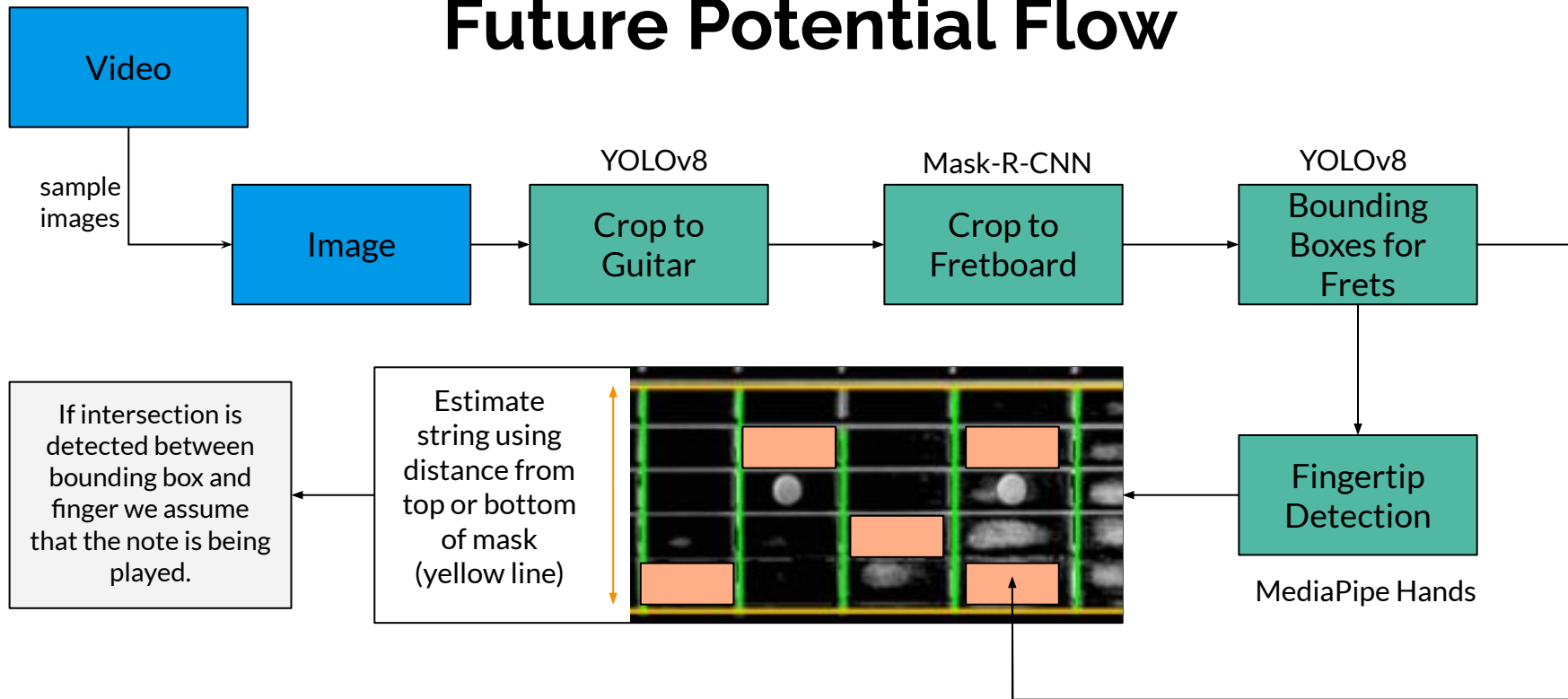
Prediction for 31_8.png



Prediction for 34_7.png



Future Potential Flow



Thank you!



Advisors: Professor Yung-Hsiang Lu and Professor Yeon Ji Yun
