

NOV 22ND 2024

Montgomery AI

Revised Project Proposal

Sean Kim / EE P 596 Computer Vision

Montgomery AI is an multi-stage system that automatically transcribe guitar tab scores from images/videos of guitar performances




PROBLEM RELEVANCE

Learning a skill alone is not easy,
especially in music. Can AI help?


More people are trying to learn instruments alone but the challenge of replicating what they see can be overwhelming. Automatic transcription can save time for learners.

OVERVIEW



INPUT

DESIRED OUTPUT



Billie's Bounce Wes Montgomery Solo, with Tablature Transcription for Guitar

Brandon Wall
1.64K subscribers

Subscribe

1.1K | Share | Download | Save


Do this automatically for ANY guitar playing video on YouTube

PRIOR ART



Text-based LLM hallucinate often when generating a score
(Some notes/numbering are correct but the order is completely off)

PRIOR ART




AnthemScore 5


Music AI for Your PC

AnthemScore is software for automatic music transcription using AI. Convert audio files like MP3 and WAV into sheet music. No subscription required—buy once and use forever on your own machine running Windows, Mac, or Linux. Free trial with 30 sec of sheet music.

[Download Trial](#)

[Buy](#)





Home Products Solutions About Us Help English



Overview Features Pricing Download FAQ LOGIN SIGN UP


Transcribe Guitar Music into TABs

Record, Upload or use a YouTube Link of Music to turn them into TABs, MIDI, MusicXML and GuitarPro files within seconds.

GET STARTED FOR FREE

Transcriptions
>4 million

Recommended by
 KEYBOARDS 

Made in Germany


14-day money-back guarantee


Existing Automatic transcription software
are mostly audio based

Why audio-based doesn't work

INPUT



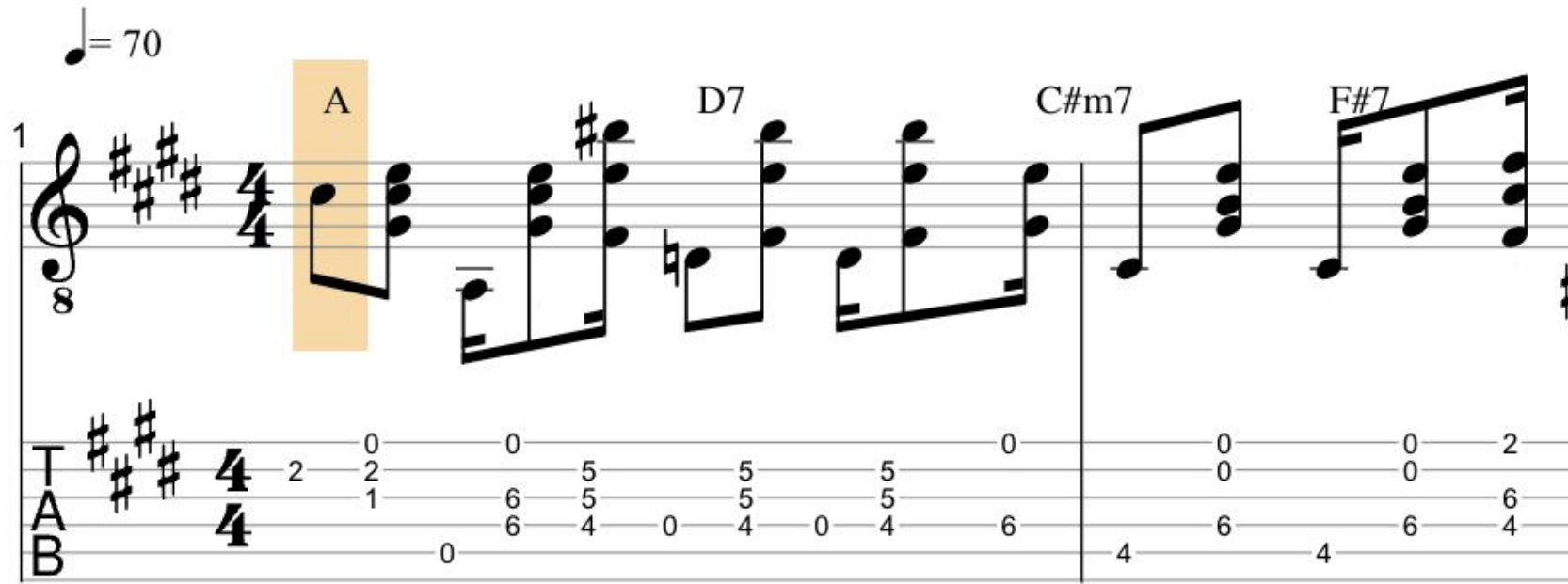
OUTPUT

 [PLAY MODE](#) [NOTE EDITOR](#) [DOWNLOAD](#)

Esperança Perdida (I Was Just One

Standard Tuning

♩ = 70



TAB

Notes are correct, but the finger position is wrong!

BASELINE

Compare the accuracy against-

1. Audio-based software: Guitar2tabs
2. Text-based LLM: ChatGPT (Ask to generate a tab score for a song)
3. Image-based LLM: LLaVa (Ask to transcribe a score given an image)

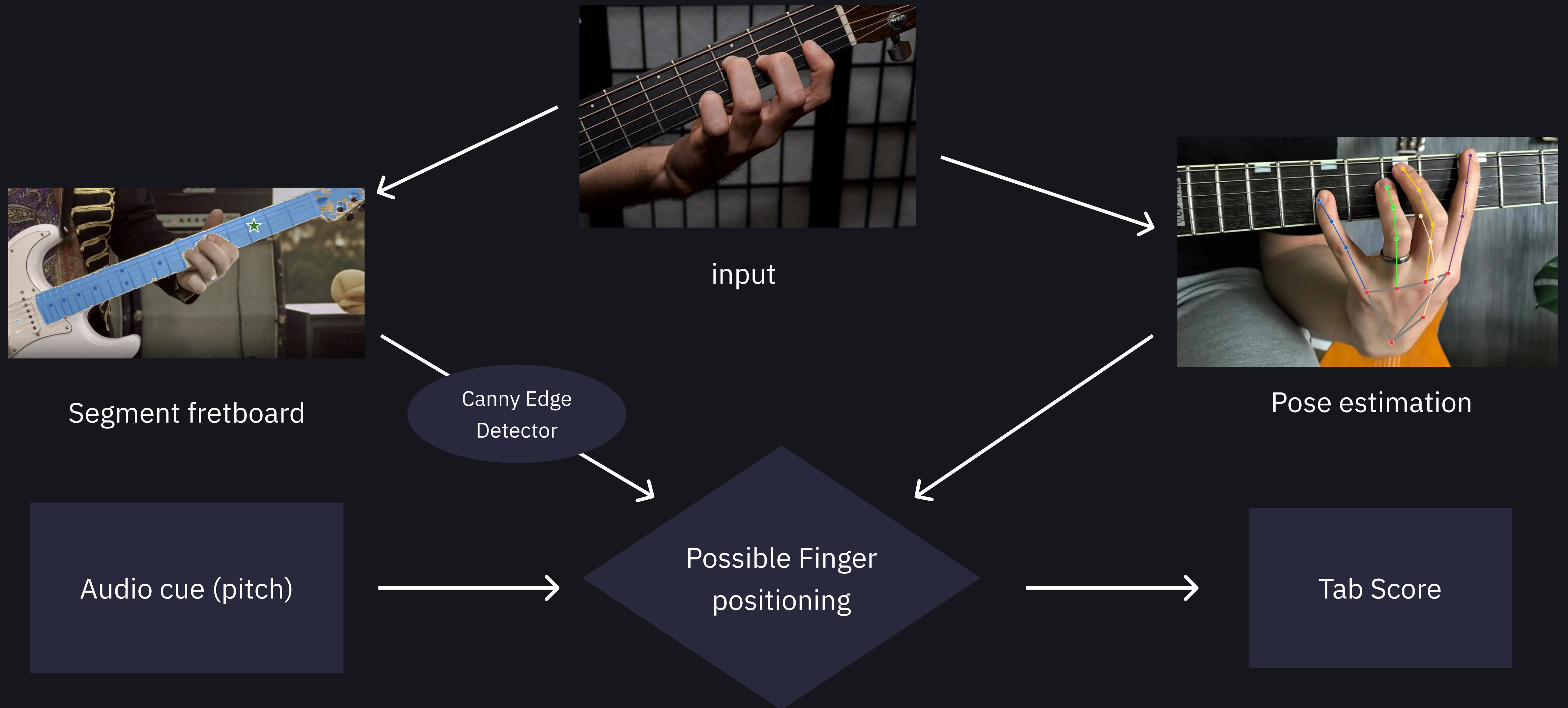
DATASET

1. Training: N/A (using pretrained)
2. Test set: Manually collect 100 images from transcribed videos on Youtube
3. (Maybe) Contact authors of relevant research papers (e.g. 3D Guitar Fingering Assessing System Based on CNN-Hand Pose Estimation...) and get the data

Proposed Methods

1. Identify fret-board using **segmentation**
2. Identify the edges inside the fret-board using **canny edge detection**
3. Identify the orientation of the fingers using **Pose Estimation**
4. Use **projective geometry** to calculate possible notes played
5. Process audio input and output a note being played on which string

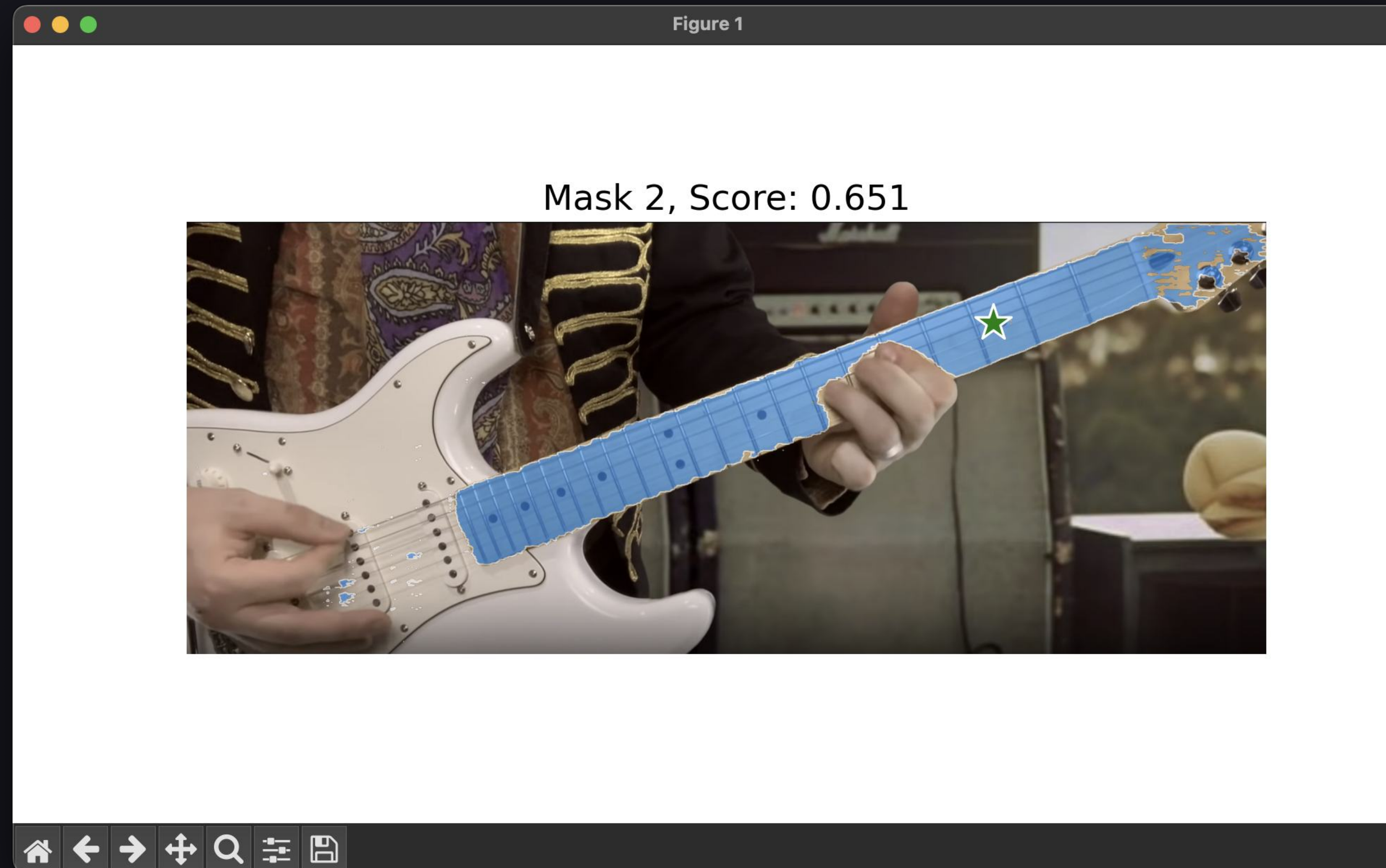
Proposed Methods (Diagram)



Milestones (by input type)

1. Image (Single Note)
2. Image (Multiple Notes)
3. Video (Single Note)
4. Video (Multiple Notes)
5. Image (Non-Orthogonal camera)
6. Video (Non-Orthogonal camera + Moving Guitar)

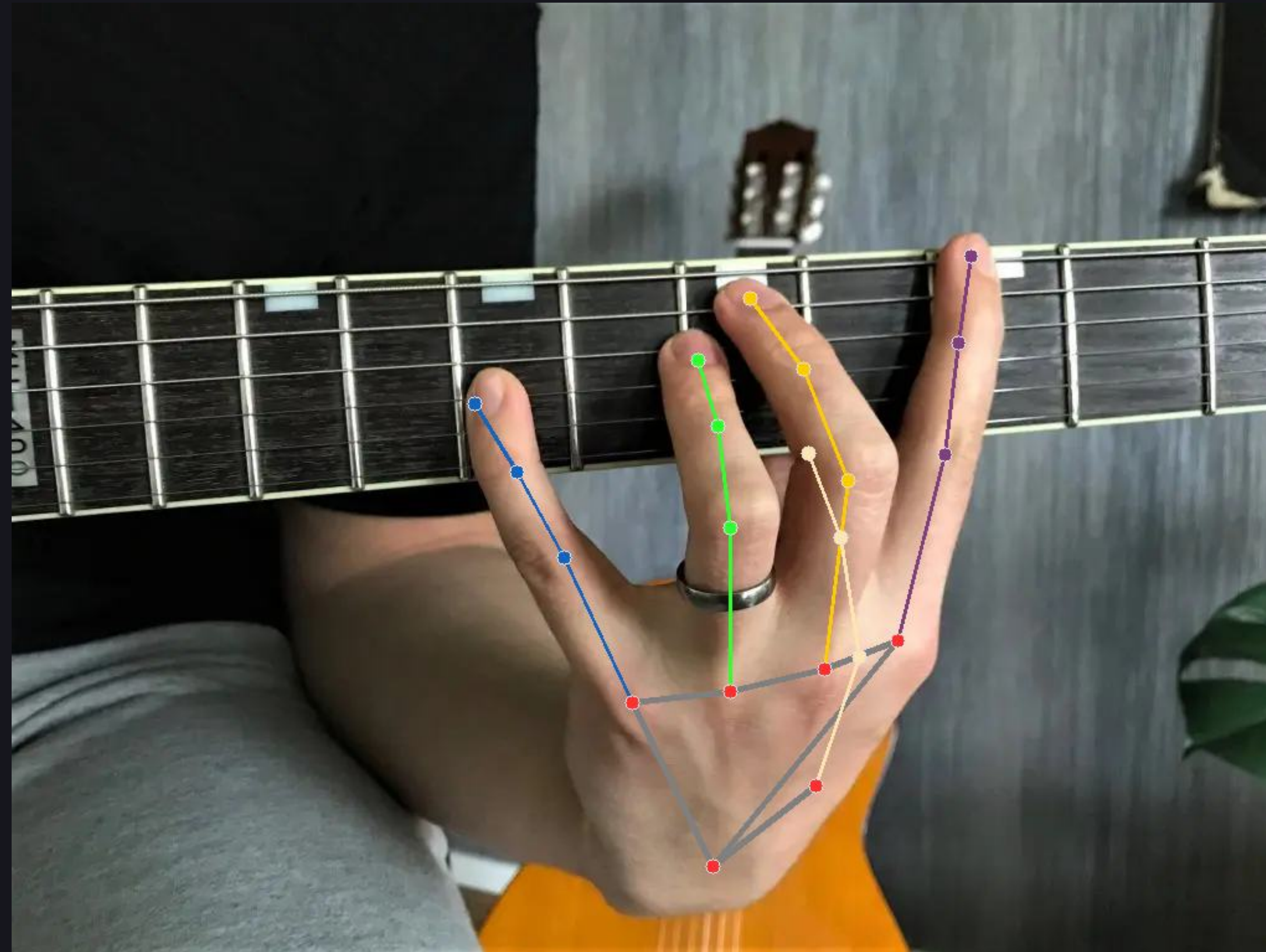
1. Fret-board Segmentation



Facebook's Segment Anything 2.1

Works reliably. Reasonable inference time on Macbook M1 Chip

3. Hand Pose-Estimation



Google's MediaPipe

Doesn't detect hand sometimes. Need to investigate

AGENDA

Future Plan

1. Refresher on problem statement
2. Update on metrics
3. Review design proposal
4. Align on GTM

Challenges / Future Plan

1. Choosing best tools: Does Segment-Anything & Media Pipe give the best result?
Is the inference time on Macbook M1 chip reasonable?
2. Choosing the correct mask: Segmentation gives multiple masks for a single point.
How can we programmatically get the “fret-board” mask?
3. Canny edge detector: would a simple canny edge detector be enough to identify the fret board bars?
4. Pose estimation customization: would training MediaPipe further on a known hand orientation (e.g. hand orientation for a playing C chord) help?
5. Audio cue: Would pitch information be enough to make decision on which fingers are actually playing the note?