

Progress Status Summary Report - FIT 3162

Team: MCS08

Project title: Singing Video Generation with Music Separation

What has been accomplished over the previous 5 weeks

In week 3, we had a discussion with our supervisor and we finally decided to use the Stylized Audio-Driven Talking-head video generation system, which is known as SadTalker, as the source code for our project. After that, we perform research on video generation models, seeking to further our understanding of the underlying concepts and techniques.

During week 4, we refine our understanding of the concepts and techniques through a structured presentation. Our supervisor provides feedback on our presentation which allows us to identify areas of improvement and solidify our knowledge. We also test various video generation models within Google Colab to explore the model capabilities and limitations.

In week 5, we convened a meeting with our supervisors to present our progress and to solicit feedback in preparation for the forthcoming interim presentation. At the same time, we discuss with our supervisor to articulate the parameters for future project enhancements, which we decide to improve the video quality, and the benchmarks for the final output.

During week 6 and 7, we were exploring and finding possible improvements for our project. We had discussions with our supervisor however some of the ways proposed do not seem feasible. Besides that, we were busy with other projects hence progress was rather slow.

In week 8, we were given an idea to improve the video quality from our supervisor. Besides that, we are now starting to create our web application UI and link it with Google Colab. As what we discuss with the supervisor, we should generate the video by using SadTalker and then get the generated video frames. After that, we need to convert each raster frame into a vector frame and then only we make those frames into video.

What degree of completion was reached in comparison to the previous set goal or plan?

Based on the Gantt Chart we created earlier, the project completion status is now almost 60% done. We have successfully done the research on our chosen facial generation model that we plan to use for our project and will be ready to perform our proposed improvements on the model. Besides that, we have started with the user interface design and will be finished by week 9 as planned. However, the improvements for video quality is still in progress as we just finished researching ways to improve the video quality, so we are going to start doing quality improvement on week 9. For further testing and debugging, we will make sure that everything is completed by the schedule written in the Gantt chart. To compare our project progress with our set goal, our expectation is much higher but we did not manage to finish it on par with our set goal. We were expected to do our prototype at week 7 but we only started it at week 8. This shows that our progress is much slower than expected which requires us to reschedule our set goal.

What was the reason for failing to meet the set goal, and what will change to avoid continuation of this situation?

There were factors affecting our progress for this project. Due to our lack of knowledge, we were having trouble finding resources on how to improve the video quality. This had held our progress for 2 weeks which is quite disappointing. We should have asked our supervisor for help earlier as he is more knowledgeable compared to us.

As a result, Our team has discussed and planned to ensure progress remains on track. In short, we had allocated more time to finish our project faster. It is also better to ask for more assistance from our supervisor so that our project might progress faster.

We are currently working on the UI and integration of the model, alongside proposed improvements on our generated video recommended by our supervisor which we discussed and agreed upon to improve our current pre-trained model, thus we believe that we are able to complete within the given time frame.

Team Member (name and ID): Toh Xi Heng (33200548)

Task attempted / completed	Completion (%)	Time taken (days or part of)	Comment, eg: reason for not completing if any
Do research for facial generation model	100%	14 days	
Test the model	100%	1 day	
Understand the model	100%	12 days	
Create user interface design	20%	1 day	That is a lot of other unit assignments the week before so we only design a simple protocol but no code it out yet.
Perform research on model improvements	80%	7 days	We have already found out the way we want to make the improvement but we need more knowledge on it before coding.
Implementing the improvements on the model	0%	0 day	As what I mentioned above, this part of work we have extended it as we are just going to done the research.

Team Member (name and ID): Yew Yee Perng (32205481)

Task attempted / completed	Completion (%)	Time taken (days or part of)	Comment, eg: reason for not completing if any
Do research for facial generation model	100%	14 days	
Test the model	100%	1 day	
Understand the model	100%	12 days	
Create user interface design	20%	1 day	Started with the implementation for UI however not done due to heavy workload
Perform research on model improvements	80%	7 days	A possible solution has been proposed and we will implement it this week.
Implementing the improvements on the model	0%	0 day	This was planned to start on week 9 thus we have yet to start but shall be completed by week 10.

Team Member (name and ID): Yeoh Ming Wei (32205449)

Task attempted / completed	Completion (%)	Time taken (days or part of)	Comment, eg: reason for not completing if any
Do research for facial generation model	100%	14 days	
Test the model	100%	1 day	
Understand the model	100%	12 days	
Create user interface design	20%	1 day	Since we just finished our research on how to improve the video quality, we have insufficient time to create a robust UI design.
Perform research on model improvements	80%	7 days	Although we understand the concept on how we should improve the video quality, we haven't found a way to integrate the model into our code.
Implementing the improvements on the model	0%	0 day	This was not done due to our user interface design is still in progress.

Evidence of progress:

Agile board:

<https://student-team-jeu1flnm.atlassian.net/jira/software/projects/KAN/boards/1?atlOrigin=eyJpIjoiY2Q2ZTlmMGQzNWZhNDQyZDk4NjYyYnN2IiwiaWVudDUOTRmMjgiLCJwIjoiaj9>

KAN board

Search

TO DO 2

- do testing and debugging
✓ KAN-8
- Finalise and produce our deliverable
✓ KAN-16

+ Create issue

IN PROGRESS 4

- Creating a UI - image and audio input
✓ KAN-1
- Creating UI - link with colab
✓ KAN-17
- do the video quality improvement
✓ KAN-9
- Creating UI - showing generated video
✓ KAN-18

DONE 7

- Doing Research for SadTalker
✓ KAN-5
- Project Summary Report
✓ KAN-19
- Project Management Report
✓ KAN-20
- do research for video quality improvement
✓ KAN-21
- Testing on video generation models
✓ KAN-22
- interim presentation slide
✓ KAN-23

Demo Slides:

<https://docs.google.com/presentation/d/16EKsLawrT-F3qfT4dXgdn0WCYy-dsBk7TxLLrUI3Jj8/edit?usp=sharing>

Singing Video Generation with Music Separation

File Edit View Insert Format Slide Arrange Tools Extensions Help

3DMM

A statistical object model separating shape from appearance variation.

In 3DMM, the #D face shape S can be

Decoupled as:

$$S = \text{AVG}(S) + \alpha U_{\text{id}} + \beta U_{\text{exp}}$$

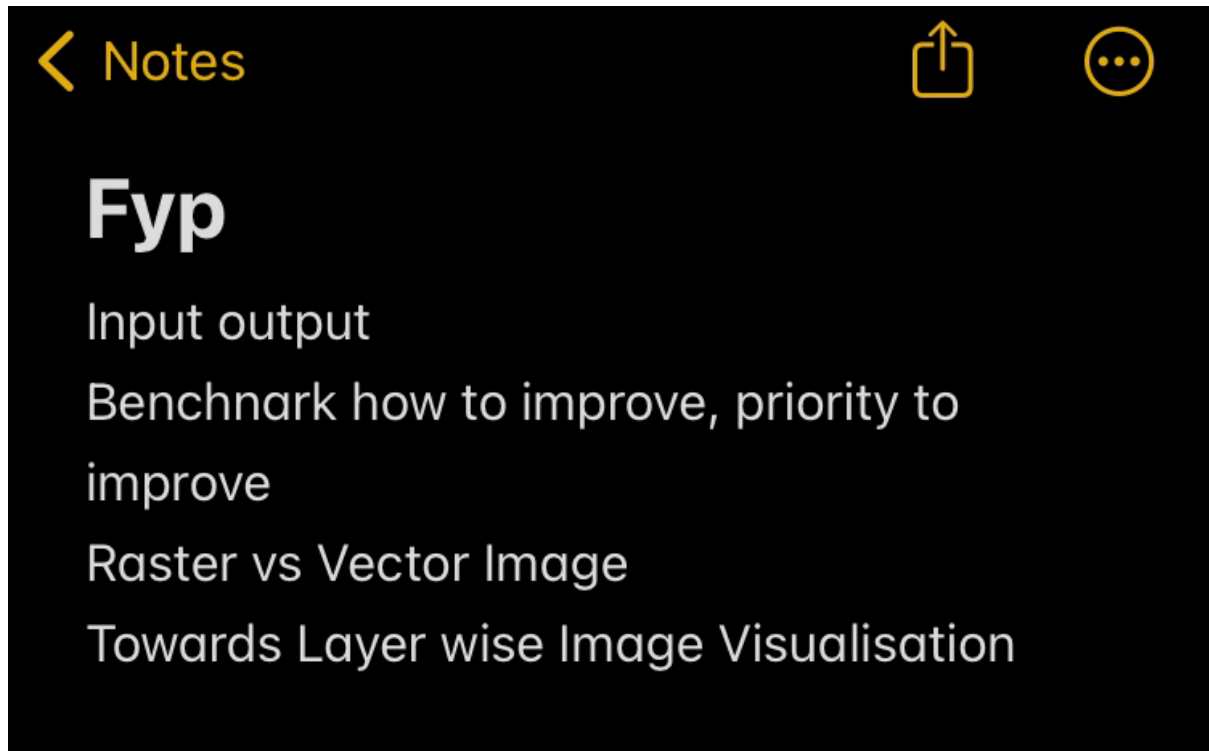
Where,

AVG(S) is the average shape of the 3D face, U_{id} and U_{exp} are the orthonormal basis of identity and expression of LSFM morphable model. Coefficients $\alpha \in \mathbb{R}^{80}$ and $\beta \in \mathbb{R}^{64}$ describe the person identity and expression, respectively.

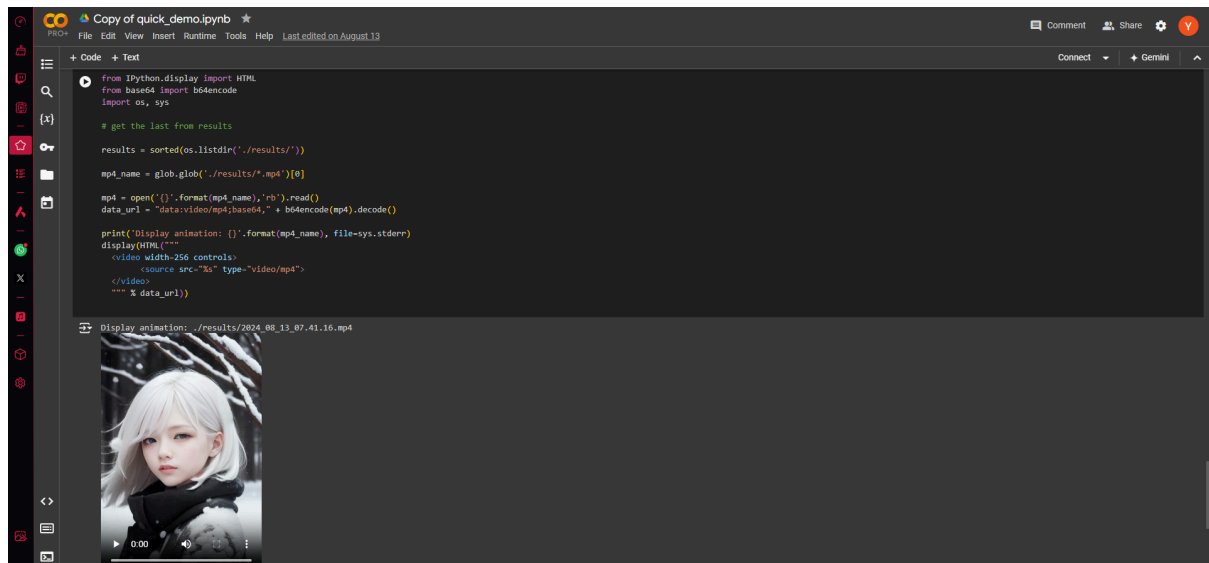
Diagram illustrating the 3DMM pipeline:

- Single Input Image (I_1, I_2) and Input Audio (d_1, a_1) are processed by Monocular 3D Face Recon.
- The output is a 3D face shape (S_1, S_2) .
- The 3D face shape is then processed by Sec. 3.2 Coeffs. Generation (Identity, Expression) to produce coefficients (α_1, α_2) and (β_1, β_2) .
- These coefficients are then processed by Sec. 3.3 3D-Aware Face Renderer to produce Generated Frames.
- Generated Frames are also processed by Expression Coefficients, Head Pose Coefficients, and Audio Feature.

Notes taken to improve video quality



Sample of generating video in colab:



Our gantt chart which indicates our current progress and our next task:

