



팀원 | B889031 방윤혁 C089067 한승호 C193256 조현석 C289004 김민현

Table of Contents

01 Background & Motivation

02 Expected Functional Output

04 System Block Diagram

06 Task Allocation for each team member

08 Target Performance Metrics

02 Example of Operational Scenario

03 Novelties of the design

05 Technical Components

07 Work Schedule

09 Performance Evaluation Method



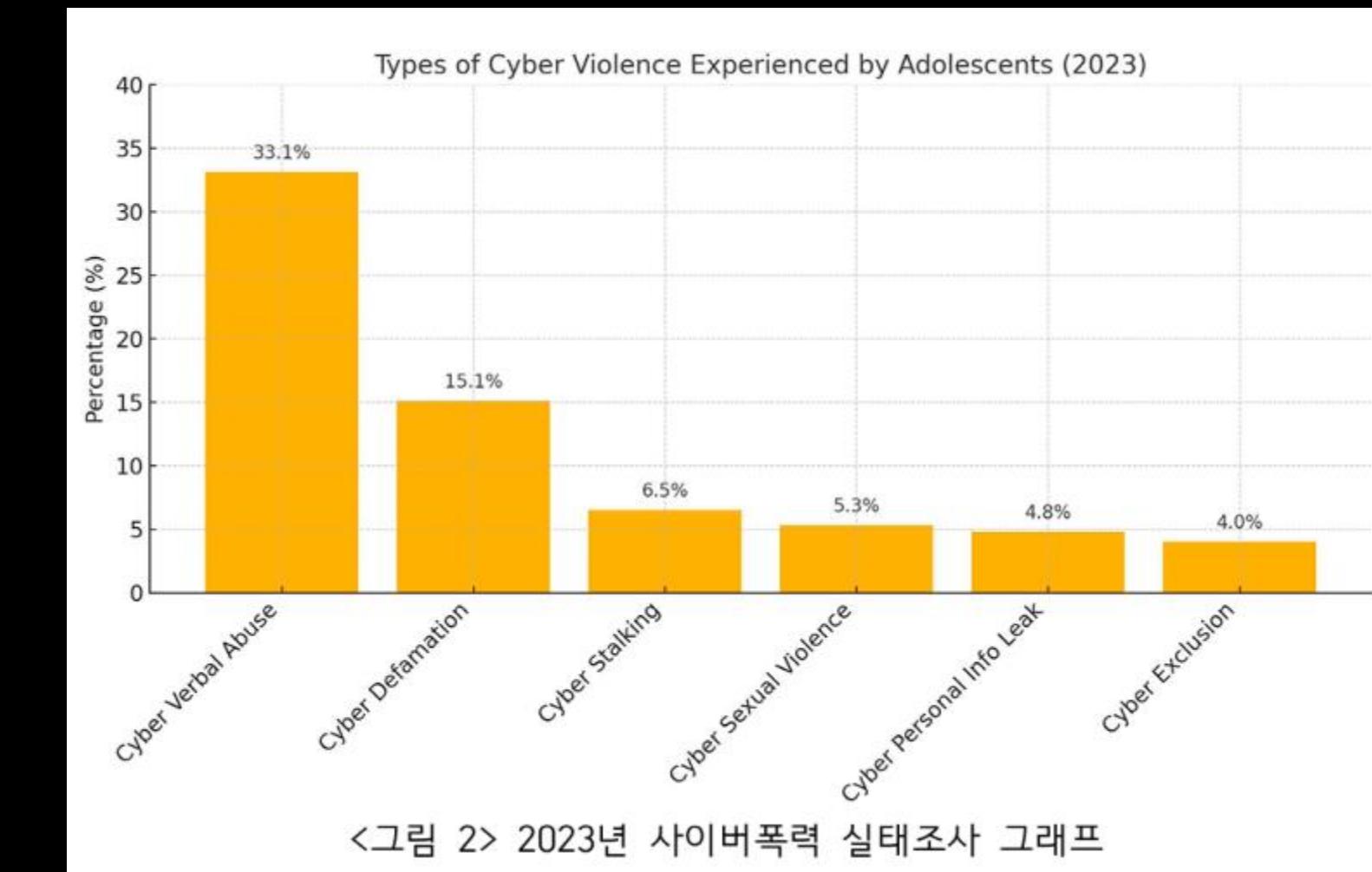
초등학생들의 희망 직업 1위는 4년째 운동선수(9.8%)가 차지하고 있다. 2위는 1년 전보다 한 계단 상승한 교사(6.5%) 였다. 3위에는 크리에이터(6.1%), 4위에는 의사(6.0%), 5위에는 경찰관/수사관(4.5%)이 각각 올랐다.

Cyberbullying and Anonymity Protection

* Key Summary

33.1% of teenagers have experienced cyber verbal abuse
(NIA 2023)

VTubers as a means of **protecting privacy**



VTubers as a way to express unique personal identities

* Key Summary

Expansion of VTuber Culture **After the Pandemic**

Enables content creation **without revealing one's face**

Utilized on metaverse platforms (e.g., ZEP, ifland)



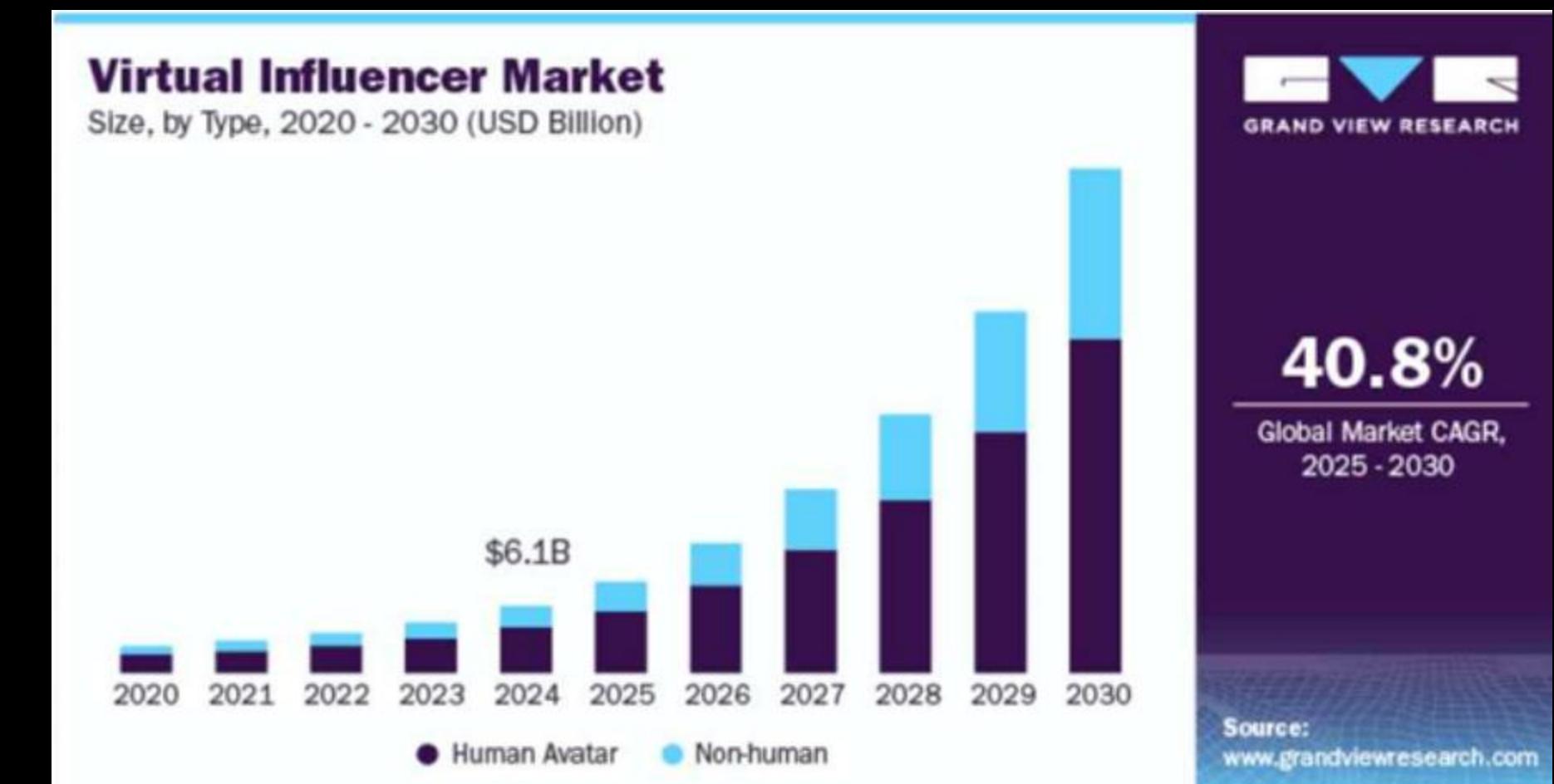
Explosive Growth of the Virtual Market

* Key Summary

2024: \$6.06 billion

→ 2030: Annual growth rate of **40.8%**

SNS Insider: Expected to reach **\$154.6 billion** by 2032



Economic Barriers to Entry for Virtual Model Creation

* Key Summary

LIVE2D commissions: Minimum of **100,000 KRW**
to over **1,000,000 KRW**

Equipment + software can cost **several million KRW**

항목	설명	예상비용
 2D 모델 제작	Live2D를 활용한 2D 아바타 제작	\$100~\$200
 3D 모델 제작	Vroid Studio 등으로 제작한 3D 아바타	\$1000~ \$5000
 하드웨어	PC, 웹캠, 마이크 등, 장비	\$500~\$2000
 소프트웨어	Live2D, Vtube Studio 등의 라이센스	\$20~\$60 /월

Virtual Model Creation



버추얼 아바타 캐릭터를 만들어줘. 강 보편적인
예쁘고 특이한 느낌으로

2D to 2.5D Conversion

There is a change!



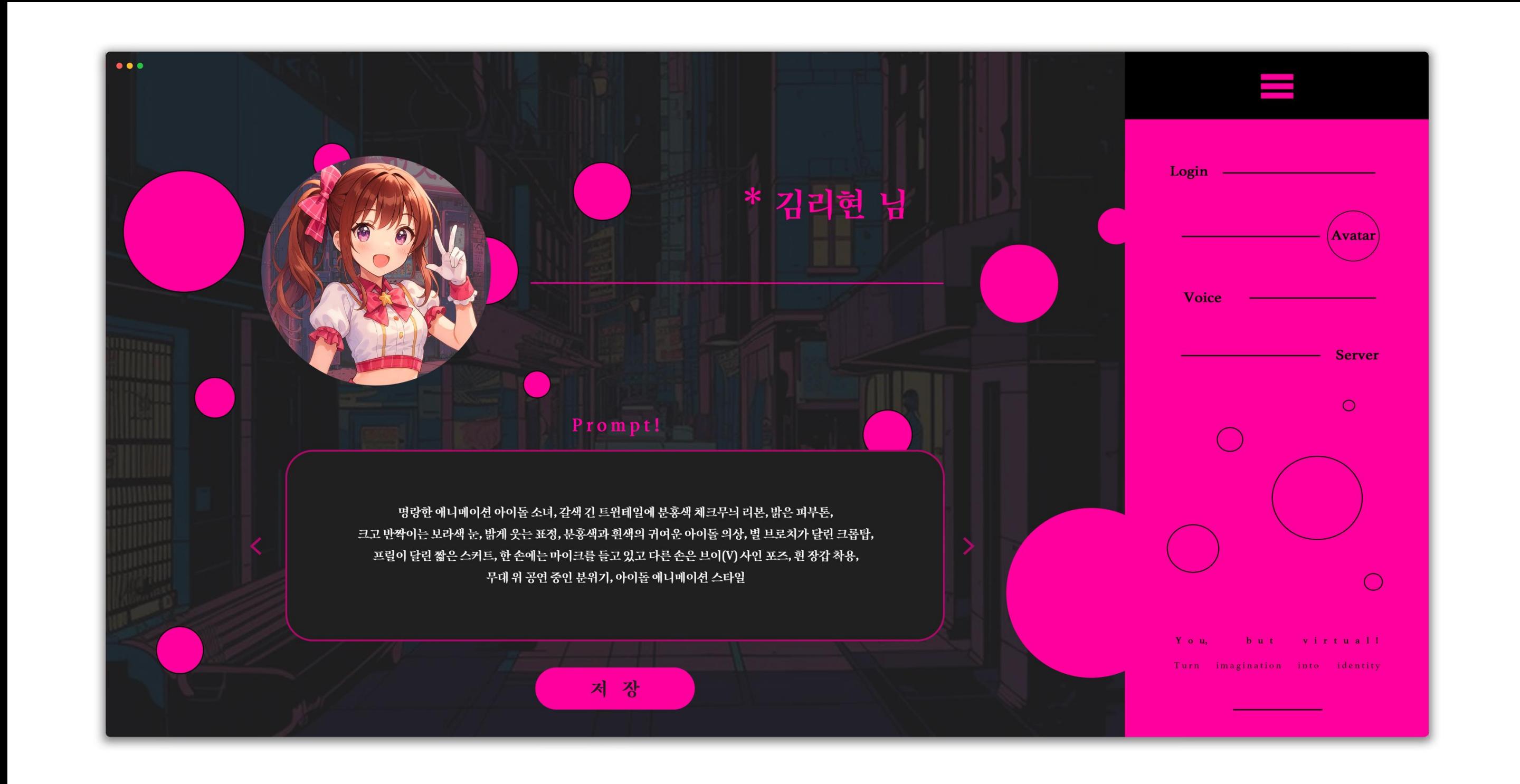
02. Example of Operational Scenario & Expected Functional Output

Final program UI example



02. Example of Operational Scenario & Expected Functional Output

Final program UI example



02. Example of Operational Scenario & Expected Functional Output

Expected Transmission Screen



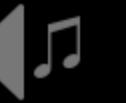
02. Example of Operational Scenario & Expected Functional Output

2D to 2.5D Conversion

There is a change!



Voice Modulation



03. Novelties of the design

Category	Virtual Model Creation	Expression-Based Model Conversion	Text Input-Based Model Creation	Voice Modification	Voice Blending	UX/UI Convenience	Expandability	Cost & Accessibility	System Integration
System (e.g., VSeeFace, FaceRig)	Only pre-made, rigged characters can be used	Only basic expressions supported, no auto changes	X	Requires external voice changer	Requires external system	Complex setup, needs additional software	Supports adding characters and modifying poses	Some software requires payment	Supports facial tracking and character export (voice modification needs external program)
iPhone Memoji/ ARKit Emoji	Fixed models only, no creation capability	Expression tracking available, but model conversion not possible	X	X	X	Simple setup, lacks expandability	Limited to Apple devices, lacks external expansion	Device-dependent, limited functionality	Supports tracking only, no export function
Auto rig pro	X	Available, semi-automated	X	X	X	Complex setup, Blender usage is complicated	Can rig other characters	Paid	Only semi-automated rigging supported
Meshy ai	Available (image to 3D model)	Available Not face rigging	Available	X	X	Easy prompt input and UI configuration	Expandable with various formats	Paid	Only generates 3D models from images

03. Novelties of the design

Category	Virtual Model Creation	Expression-Based Model Conversion	Text Input-Based Model Creation	Voice Modification	Voice Blending	UX/UI Convenience	Expandability	Cost & Accessibility	System Integration
Applio	X	X	X	Supported, but no real-time output	Voice blending supported	Limited convenience	Real-time output not supported	Free, but only offers blending feature	Voice only, no video
*Voice Changer (e.g., w-okada, used standalone)	X	X	X	Voice modification only	Voice blending supported	Limited convenience	No AR support	Free, but only offers voice modification	Voice only, no video
Live2D	2D virtual model creation specialized	Expression-based animation supported	X	No built-in voice features	No voice features	Professional UI with dedicated software	Extensive plugin support, OBS integration	Paid software (Cubism PRO required)	2D animation + streaming Integration
제안 시스템 (도전 버튜버!)	Based on LLM (ChatGPT/DALL-E), auto image generation with AR	Real-time BlendShapes → Emoji conversion	Text input → LLM → Auto image generation	Real-time AI-based voice transformation	Voice blending supported	Intuitive UI for settings	Custom prompts, supports various virtual model formats, 2D masking available	Free open-source model use, API planned	Face tracking + image generation + voice modulation + export

VTuber System Performance Analysis

01. Key Performance Indicators (Realistic)

COST REDUCTION

30-50%

\$2,000 → \$1,200

Compared to mid-tier solutions

SETUP TO DEPLOYMENT

82.5%

40hrs → 7hrs

Complete workflow

SOFTWARE REQUIRED

67%

3 → 1

Integrated solution

LEARNING CURVE

70%

20hrs → 6hrs

To basic proficiency

RESOURCE EFFICIENCY

CPU 50%

vs 60-80%

Lower usage vs standard

COST SAVINGS ROI

400%

4x investment

Hardware ownership assumed

02. Cost Comparison by Quality Tier

Quality Tier	Traditional Market Price	Our System Cost	Realistic Savings
Basic Entry Level Simple model, basic features	\$500-800 (Commission \$300 + Tools \$200)	\$0 (Only PC + iPhone required)	100% (If you have hardware)
Mid-Tier Professional Custom design, good rigging	\$1,500-2,500 (Commission \$1,000 + Rigging \$500 + Tools)	\$0 (Only PC + iPhone required)	100% (If you have hardware)
High-End Premium Professional quality, complex rigging	\$3,000-5,000 (High-end commission + expert rigging)	\$0 (Only PC + iPhone required)	100% (If you have hardware)

💡 **Single Configuration:** Our system offers one unified setup that works for all quality tiers. No need for expensive upgrades or additional software. If you already own a PC (mid-spec) and iPhone (Face ID enabled), you can start immediately at zero additional cost.

03. Real-World Streaming Performance Analysis

Based on Actual User Reports (2025): Analysis of VTube Studio, VSeeFace, and other popular solutions with 100+ hours of community feedback and testing data.

Platform	Real Performance	Common Issues	Our System Comparison
VTube Studio (Most Popular)	<ul style="list-style-type: none">• 30-45fps typical• CPU: 35-50%• RAM: 3-4GB• Easy setup	<ul style="list-style-type: none">• Frame drops during gaming• Lag with high-spec games• Expression tracking limited• CPU bottleneck common	Similar: 30fps, CPU 50% Better: More stable Advantage: Zero cost if have hardware
VSeeFace (Advanced Users)	<ul style="list-style-type: none">• 45-60fps capable• CPU: 40-60%• RAM: 4-5GB• VRM compatible	<ul style="list-style-type: none">• Complex setup• Audio-video sync issues• Driver instability• High hardware demand	Better: Simpler setup Better: Zero cost Worse: Lower FPS currently
Live2D + nizima (Professional)	<ul style="list-style-type: none">• 60fps smooth• CPU: 30-45%• High quality• Expensive	<ul style="list-style-type: none">• \$1,500-3,000 initial cost• Steep learning curve• Time-consuming setup• Multiple software needed	Better: 100% cost savings Better: Faster setup (7hrs) Trade-off: Lower FPS for now
AI Motion Capture (Viggle LIVE, etc.)	<ul style="list-style-type: none">• 30-45fps• AI-powered• Quick start• Limited customization	<ul style="list-style-type: none">• Precision limitations• High system requirements• Less customization• Subscription costs	Better: More customization Better: No subscription Similar: AI approach

04. Real User Scenarios & Solutions

Streaming Scenario	Common Problems (Industry)	Our System Performance	Optimization Tips
Gaming + VTubing (AAA titles)	<ul style="list-style-type: none">Frame drops to 15-25fpsCPU: 80-95% usageFrequent stutteringRequires high-end PC	<ul style="list-style-type: none">Maintains 25-30fpsCPU: 70-85% totalSmoother experienceMid-tier PC sufficient	<ul style="list-style-type: none">Lower game graphics to HighUse 720p streamingClose background apps
Chatting Stream (Low CPU games)	<ul style="list-style-type: none">Generally stableCPU: 35-50%No major issuesAll solutions work well	<ul style="list-style-type: none">Very stableCPU: 30-40%Excellent multitaskingCan run Discord, browser, etc.	<ul style="list-style-type: none">Ideal use casePlenty of headroomAdd overlays freely
Singing/Music (Audio quality critical)	<ul style="list-style-type: none">Audio-video sync issuesBuffer adjustments neededLatency problemsComplex audio routing	<ul style="list-style-type: none">140-180ms latencyAcceptable for singingRequires audio interfaceManual calibration needed	<ul style="list-style-type: none">Use ASIO driversAdjust OBS audio delayTest before streamMonitor audio closely
Drawing/Art (Screen sharing)	<ul style="list-style-type: none">High CPU loadCanvas capture issuesDelay in movements45-65% CPU typical	<ul style="list-style-type: none">Lower CPU overheadSmoother canvas capture40-55% CPU usageBetter responsiveness	<ul style="list-style-type: none">Use window captureReduce canvas sizeLower streaming bitrate
IRL/Outdoor (Mobile setup)	<ul style="list-style-type: none">Limited to iPhone/iPadBattery drain severeNetwork instabilityQuality compromises	<ul style="list-style-type: none">iPhone Face ID worksBattery: 2-3 hours max4G/5G streaming possibleLower quality but functional	<ul style="list-style-type: none">Bring power bank (20,000mAh+)Test network beforehandUse 720p 30fpsHave backup plan

20장 부연설명

Common Streaming Issues & Solutions

- ⚠ **Issue:** Face tracking freezes during high CPU spikes → **Solution:** Set VTuber software to "High Priority" in Task Manager
- ⚠ **Issue:** Audio out of sync after 30+ minutes → **Solution:** Restart stream or adjust OBS audio delay (+50-100ms typically needed)
- ⚠ **Issue:** Expression recognition poor in dim lighting → **Solution:** Add ring light or desk lamp, ensure face is well-lit
- ⚠ **Issue:** Frame drops when switching scenes in OBS → **Solution:** Pre-load scenes, reduce transition effects, use Studio Mode
- ⚠ **Issue:** Memory leak after 3+ hours streaming → **Solution:** Schedule 5-min break every 2 hours, restart software
- ⚠ **Issue:** Overheating laptop/PC during summer → **Solution:** Laptop cooling pad, room AC, reduce CPU load (lower game settings)

05. Feature Comparison (Honest Assessment)

Feature	Traditional Solution	Our System	Assessment
Character Generation	2-7 days Artist commission	5 minutes AI-generated	99% faster Instant creation
Rigging Process	8-40 hours Manual work	6 hours Semi-automated	80-85% faster Good quality
Face Tracking FPS	45-60fps Industry standard	30fps iPhone-based	Below standard Roadmap: 45fps+ by Q3 2025
CPU Usage	60-80% During gaming streams	50% More headroom	Better multitasking Gaming-friendly
Setup Complexity	High Multiple software	Medium Integrated solution	Easier but still requires 6-hour learning
Customization Level	High Full artistic control	Medium AI-assisted	AI limits some artistic nuances

06. AI Automation Performance (Conservative Estimates)

Process Stage	Traditional Time	Our System Time	Automation Level	Quality Note
Character Design	24-72 hours	5 minutes	99.4% faster	Instant generation
Initial Rigging	8-20 hours	Instant (automated)	100% automated	Auto-generated
Expression Setup	4-10 hours	Instant (automated)	100% automated	Good default quality
Voice Training	2-4 hours	3 hours	Requires training time	Quality depends on input
Testing & Integration	8-16 hours	6 hours	50-62% faster	Testing required
Total End-to-End	46-122 hours	~7 hours	85-94% faster	Production-ready

07. Technical Performance (Measured Results)

Metric	Target	Current	Industry Standard	Status
Face Tracking FPS	60fps	30fps	45-60fps	Needs Improvement Roadmap: Q3 2025
CPU Usage	<50%	50%	60-80%	Better efficiency
Memory Usage	<3.5GB	2.5-3.5GB	4-6GB	30-40% Better
GPU Usage	<40%	28-38%	45-65%	25-35% Better
Pipeline Latency	<120ms	140-180ms	150-200ms	Acceptable Room for improvement

💡 Key Strength: Our system excels at resource efficiency and automation speed. Character generation takes 5 minutes instead of days, and the entire setup process is completed in ~7 hours instead of weeks.

08. Cost Savings Analysis (Not Revenue)

Scenario	Traditional Cost	Our System Cost	Savings	Notes
With Hardware (PC + iPhone owned)	\$1,500-2,500	\$0	100% (\$1,500-2,500)	Most realistic scenario Zero additional cost
Need iPhone	\$2,200-3,200	\$700	68-78% (\$1,500-2,500)	iPhone Face ID required
Complete Fresh Start	\$3,200-4,200	\$1,700	47-60% (\$1,500-2,500)	PC + iPhone purchase needed

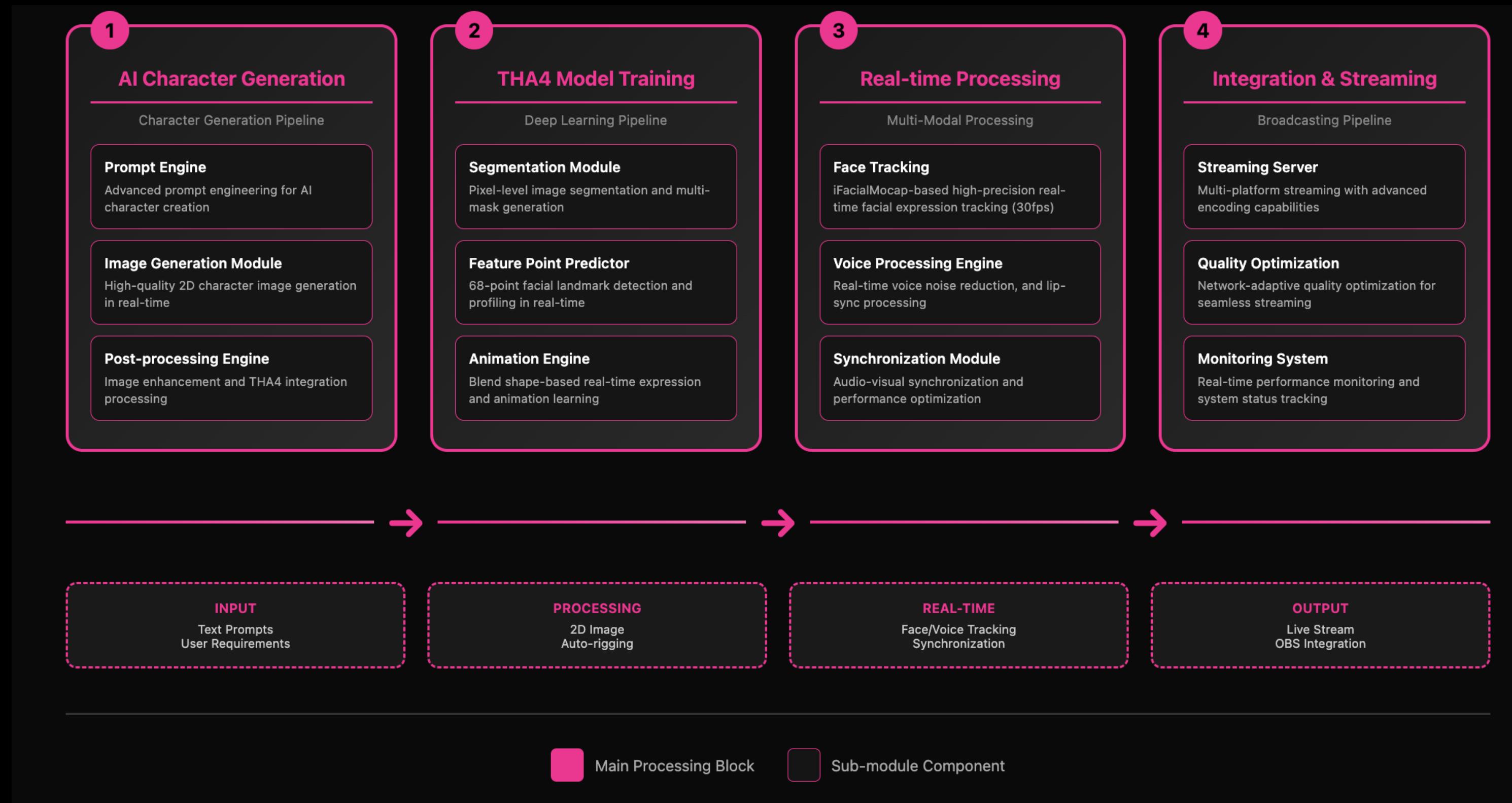
⚠ Important Clarification: The percentages above represent **cost savings ONLY**, not revenue generation. VTuber monetization is challenging:

- YouTube Partner Program requires 500 subscribers + 3,000 watch hours
- Average ad revenue: \$0.05-0.70 per 1,000 views (extremely low)
- Super Chat income highly variable and audience-dependent
- Most VTubers take 6-18 months to achieve basic monetization

This system reduces initial investment costs but does NOT guarantee income.

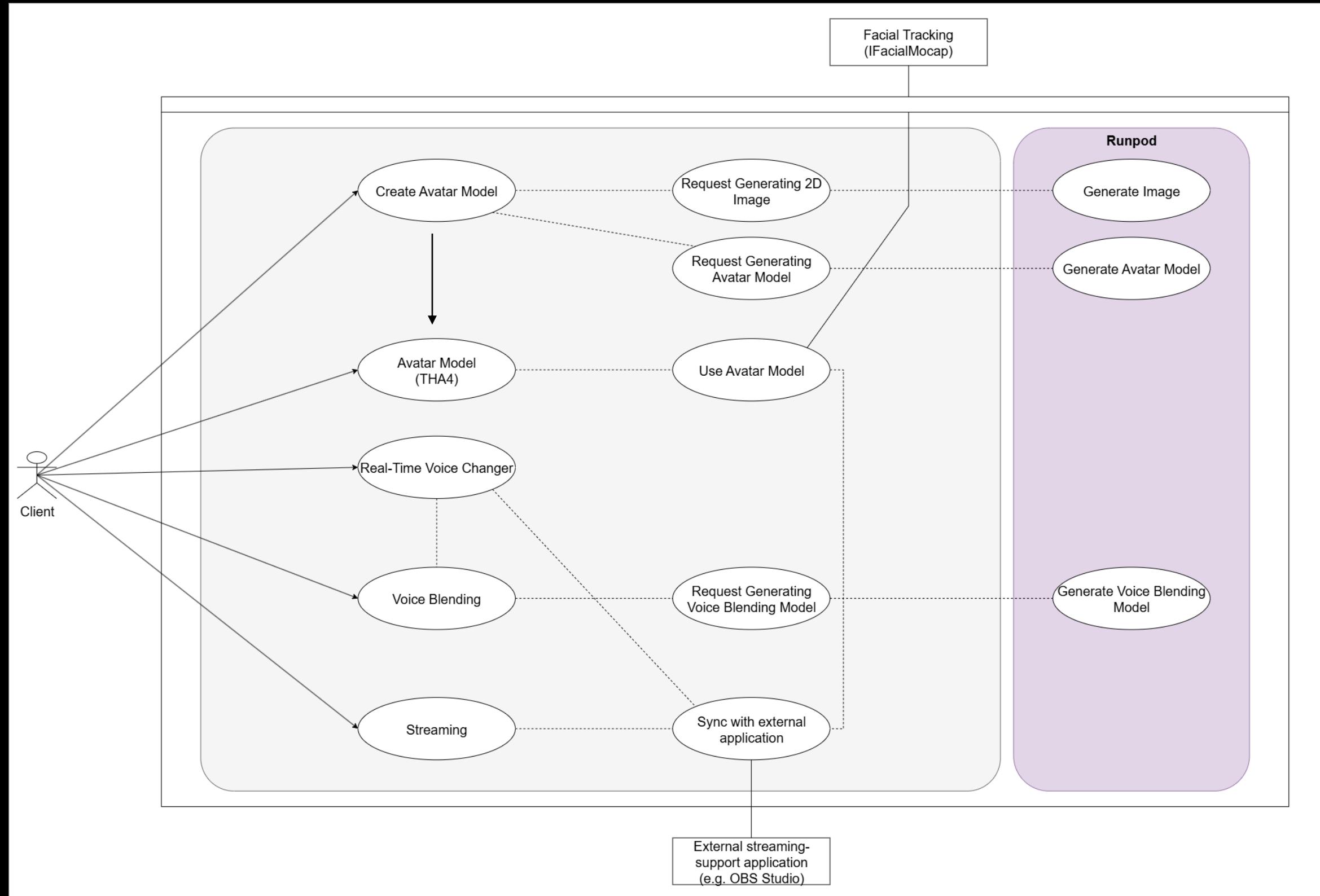
System Block Diagram

There is a change!



System Block Diagram

There is a change!



Virtual Voice

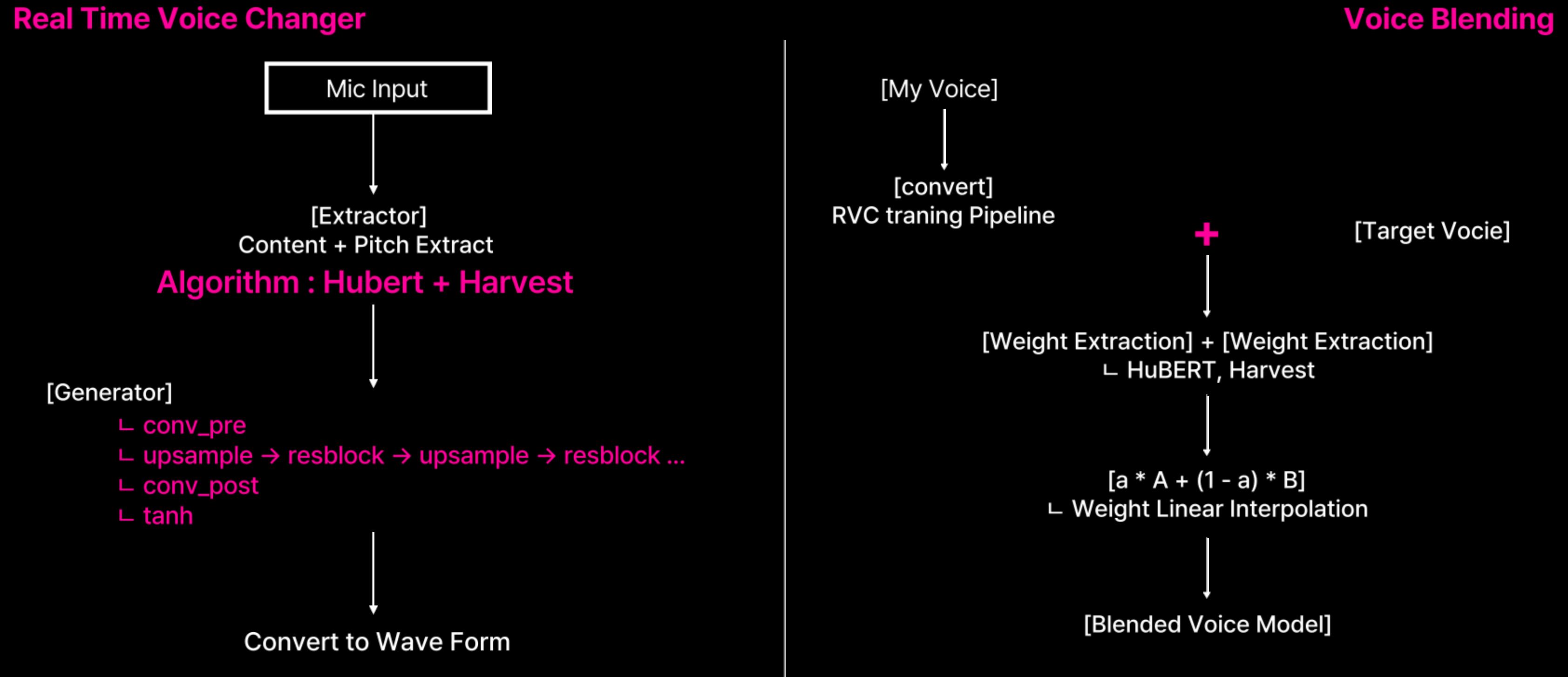


Image Generation

Prompts

1girl, hatsune miku, worst quality, ms paint (medium).



1girl, happy, smile, crying, 2koma.

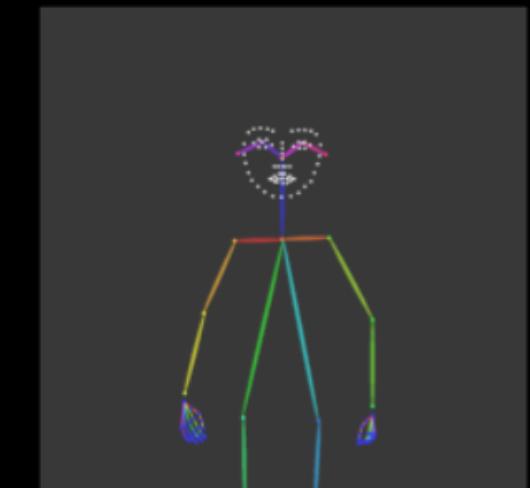
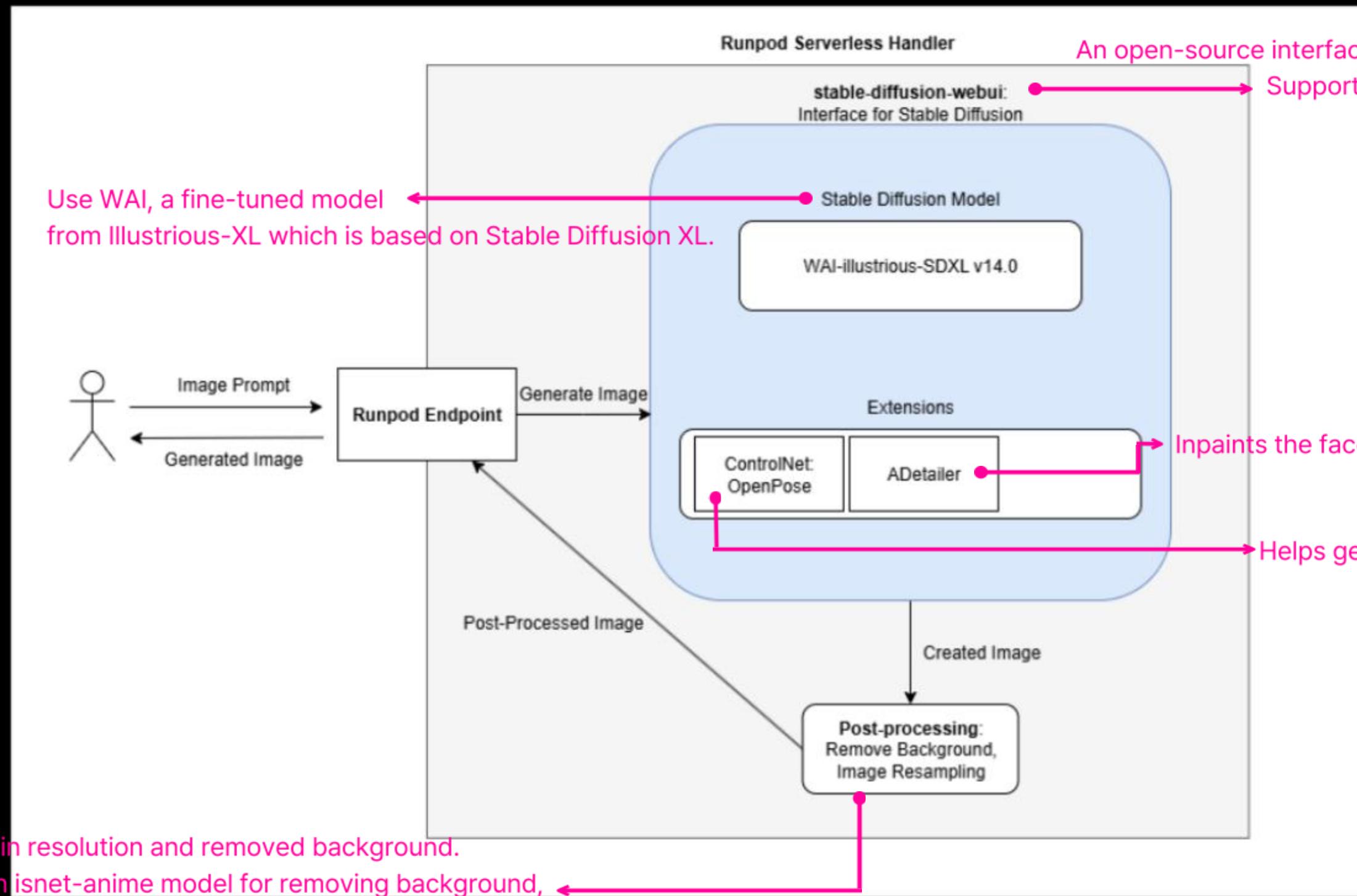


Fine Tuning (LoRA)



Image Generation

There is a change!



To generate a model with THA4, an input image should have certain resolution and removed background. We use rembg Python library with isnet-anime model for removing background, and Pillow for image resampling.

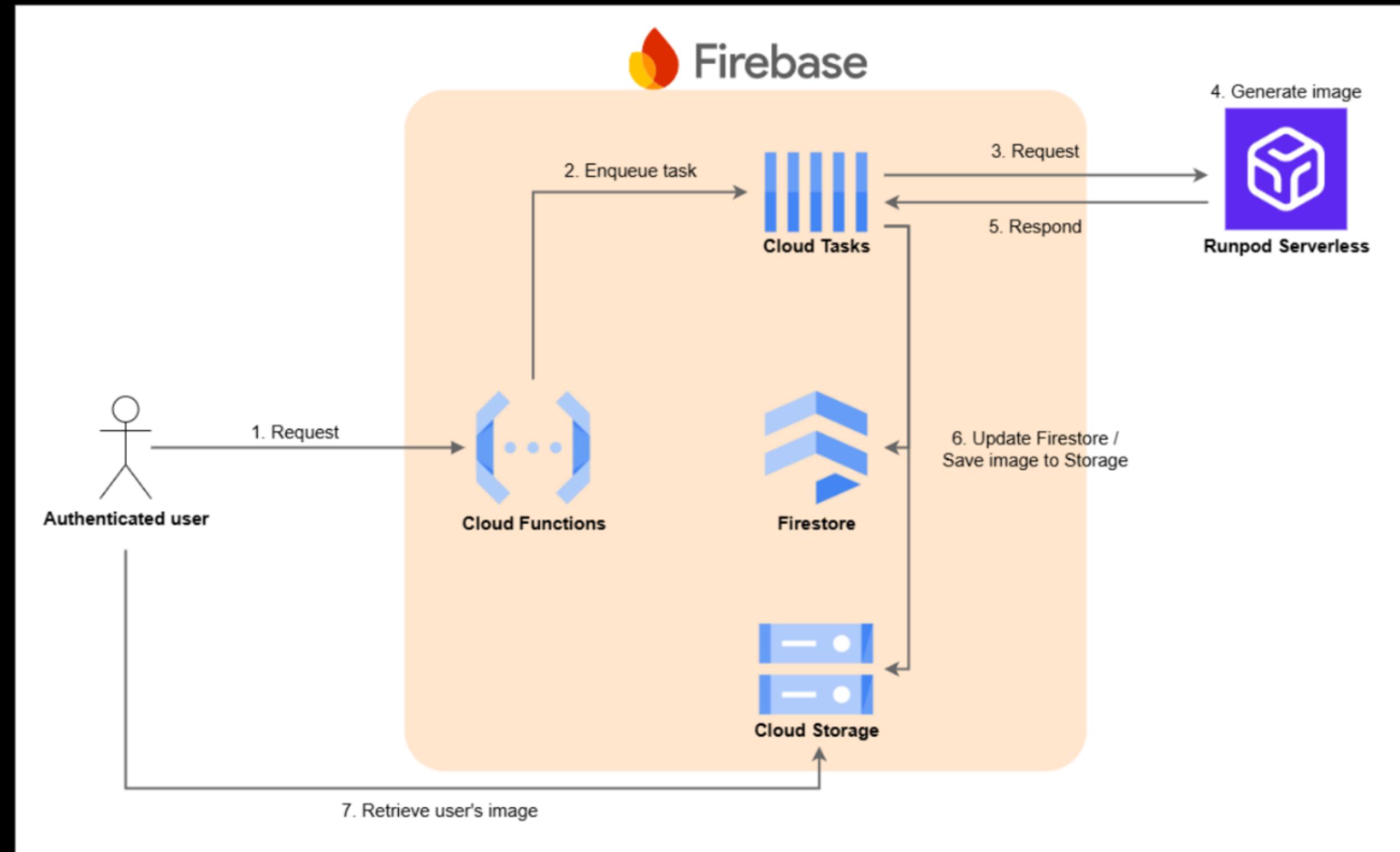
Image Generation

There is a change!



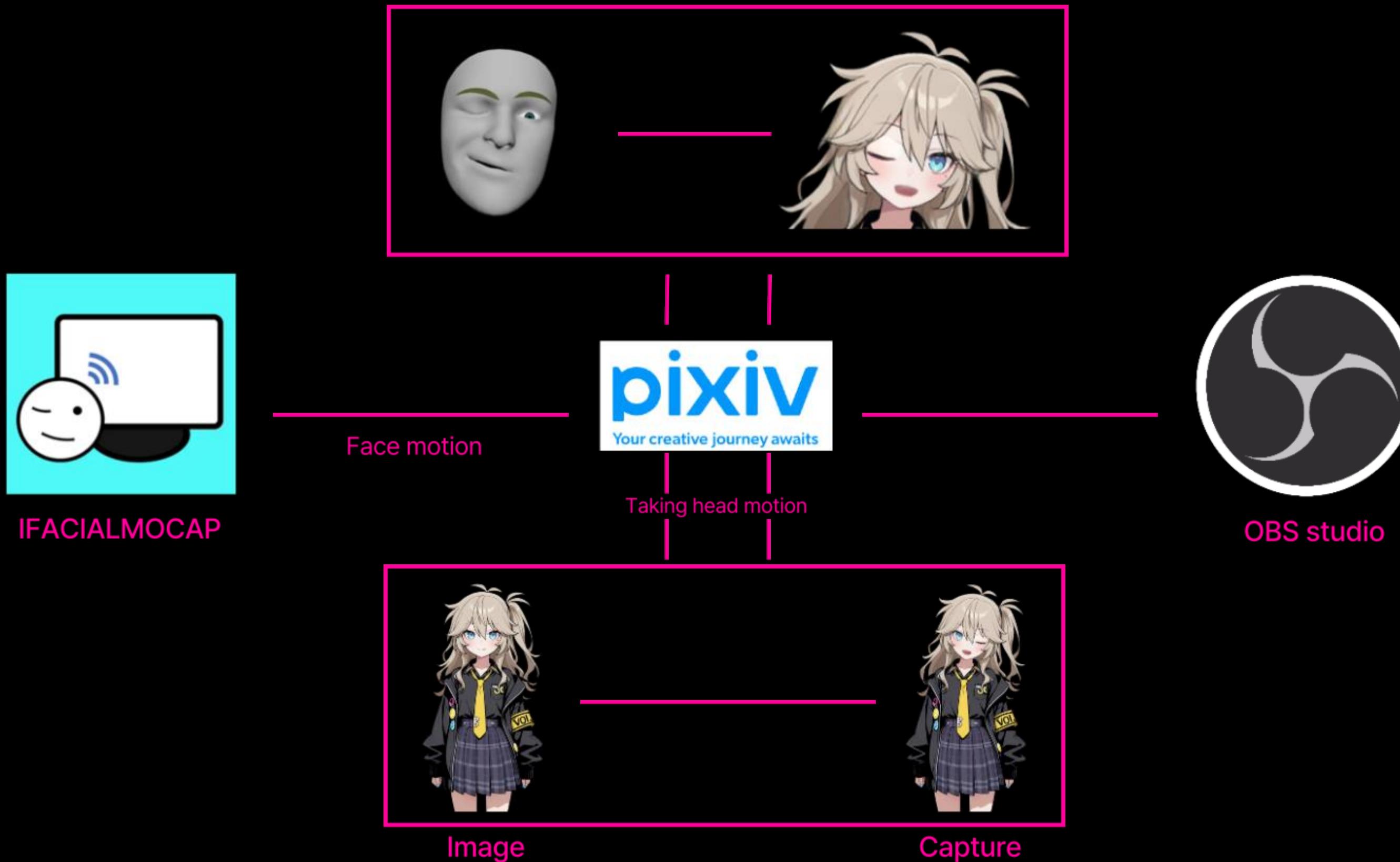
Image Generation

There is a change!



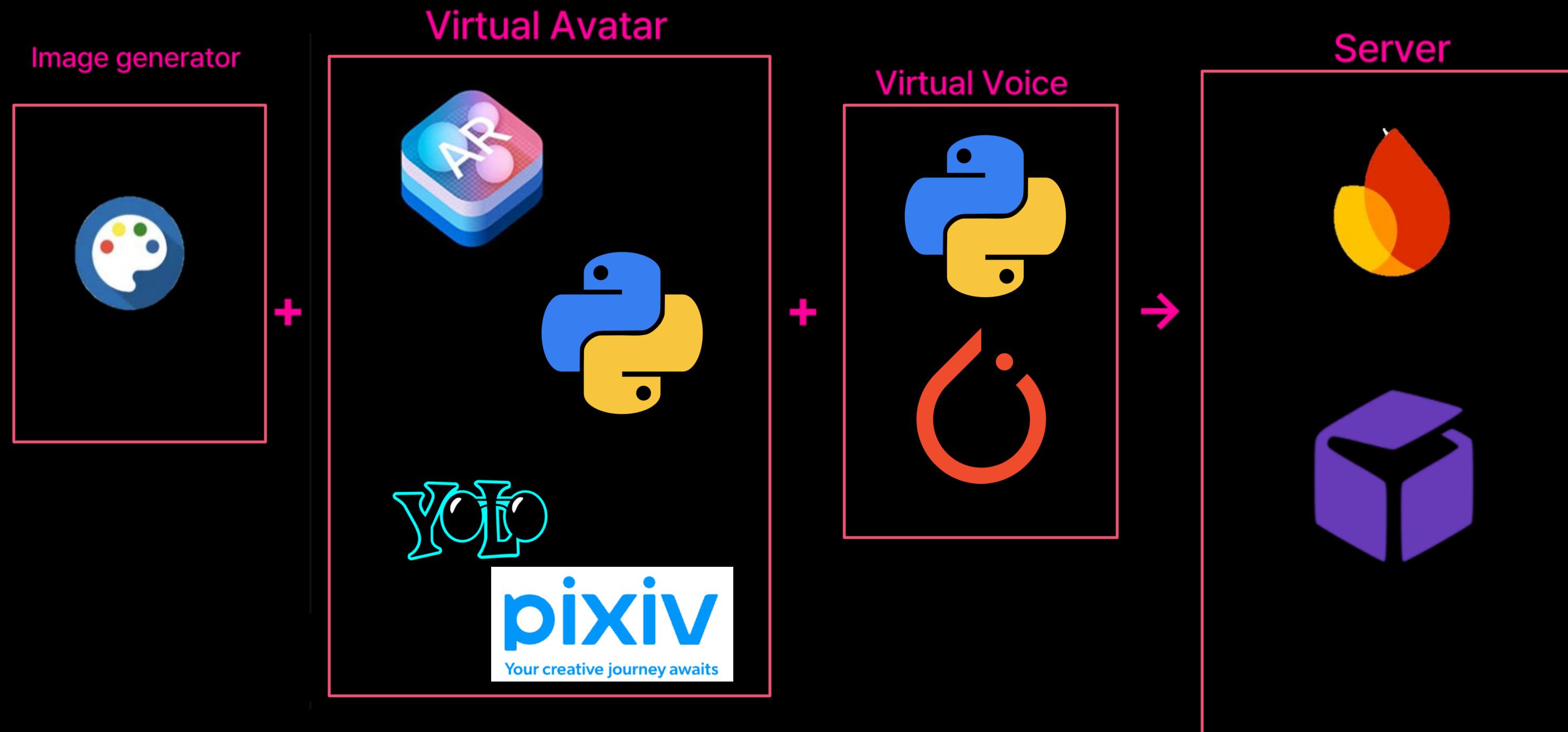
Modeling

There is a change!



Virtual Avatar

There is a change!



Task Allocation for each team member

01 방문혁

AR integration, automated image rigging,
iPhone compatibility

02 한승호

Image rigging model tuning,
Server implementation

03 조현석

Voice modulation,
Server implementation

04 김리현

UI development,
LLM model training

Work Schedule

Task Description	May	Jun	Jul	Aug	Sep	Oct	Nov
Image and model generation based on user request							
Real-time voice changer and voice blending implementation							
Development of auto-rigging system							
Backend system construction							
Development of user application							

Target Performance Metrics

Functions	Detailed Items
User APP	- Animates model based on face tracking and displays the result in the screen
Real-Time Voice Changer	- Performs real-time voice modification using voice feature mapping, with noise filtering
Voice Blending	- Blends the user's voice with a target voice based on tone, without altering pitch
Creating VTuber Model	- Generates a character image from LLM-processed text prompts, and creates a model from the generated image
Auto-Rigging System	- Automatically performs facial rigging by segmenting facial features and attaching rigging points to each feature
Server	- Loads AI models and the database, then runs inference on the models based on user's inputs

Performance Evaluation Method

Functions	Detailed Items
Real-Time Voice Changer	<ul style="list-style-type: none">- Latency: Voice delay from input to output. Should be under 100ms for smooth real-time use.- PESQ: Checks how natural the voice sounds after conversion. Over 3.5 means good quality.- STOI: Measures how clearly the voice can be understood. Over 0.75 is considered clear.
Voice Blending	<ul style="list-style-type: none">- Output Stability: We check if voice quality stays stable with different input voices. If scores don't change much, the model is considered stable.
VTuber Image Generation	<ul style="list-style-type: none">- CLIP Score: Measures how well the generated image matches the user's text prompt (like eye color, hair style). Over 0.30 is a good match.
Auto Rigging	<ul style="list-style-type: none">- We compare our face landmark model with other public models to see which works better for animation.