

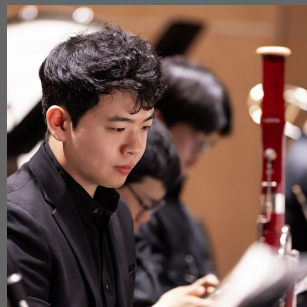
About Supertone



Neural analysis & synthesis
Unified voice synthesis framework

About Speaker

- Previously Game Engine Programmer. Used to develop GPU Graphics pipeline or Physics Engine with C++.
- Joined Supertone at 2022. Developer of User-centric products.
- Currently the Leader of the '**AudioDev team**'





The Voice Intelligence Platform

Application
Web / Local

uc
CLEAR
VOICE SEPARATOR
Real-Time Voice Separator for Voice Separation

SHIFT
REAL-TIME VOICE CHANGER
Real-Time Voice Changer for Interactive Content Creation

Project
SCREENPLAY
Text to Speech

Project
VOICE MARKET PLACE

VOICE LIBRARY

Solution

API

TTS

CVC

SVS

SDK

RTVC

RTSE

Framework

NANSY
Neural analysis & synthesis
Unified voice synthesis framework

NAUD
Efficient development of voice AI-based applications



Introduction

Quick Demo

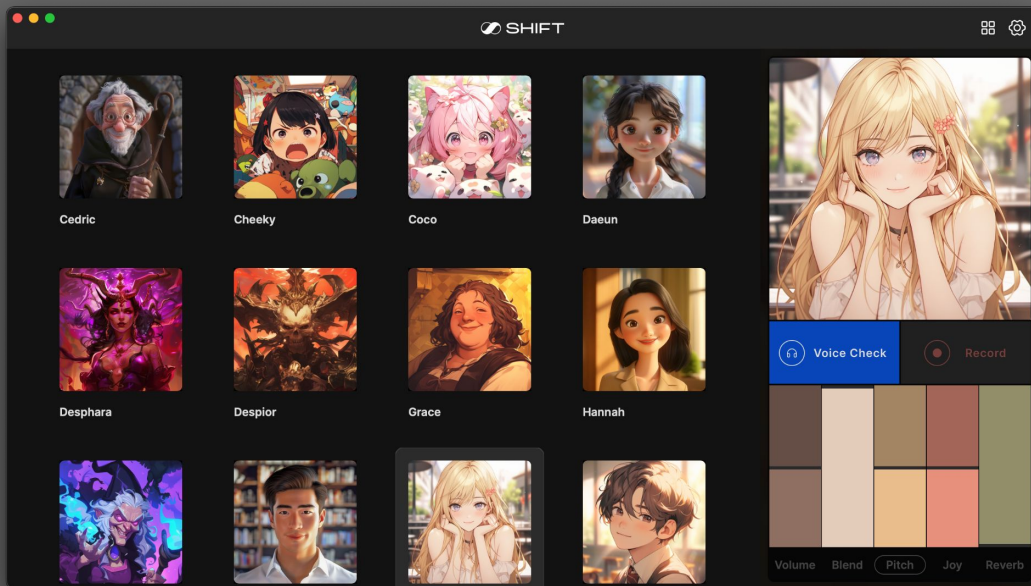
Supertone Clear

Real-time voice separator Plug-in



Supertone Shift

Real-time voice conversion Software




Supertone Air

Reverb & EQ Dialogue Match





 SUPERTONE

AIR REVERB & EQ
DIALOGUE
MATCH

Capture the Air In Dialogue





Deep Dive

Supertone Clear

Real-time Speech Enhancement

Data preprocessing to obtain clean voice

Recorded Voice
with Noise

RTSE

NANSY

Neural analysis & synthesis

Unified voice synthesis framework

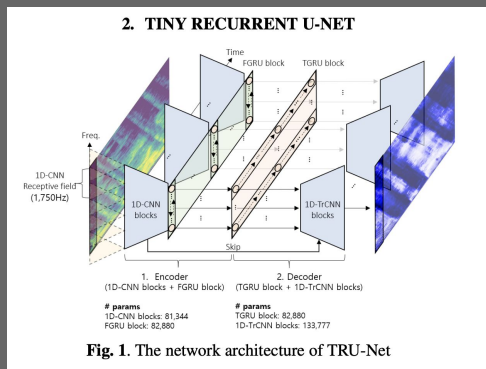


Fig. 1. The network architecture of TRU-Net

REAL-TIME DENOISING AND DEREVERBERATION WITH TINY RECURRENT U-NET *H, Choi et al*

- Lightweight enough to run in real-time
- Effectively separates non-voice and wet reverberant voice.
- Great starter to experiencing a Product development.

Developing a Plug-in

Definitely Fun

Drawing a prototype
Experimental Features
Pricing a plug-in

Mostly Fun

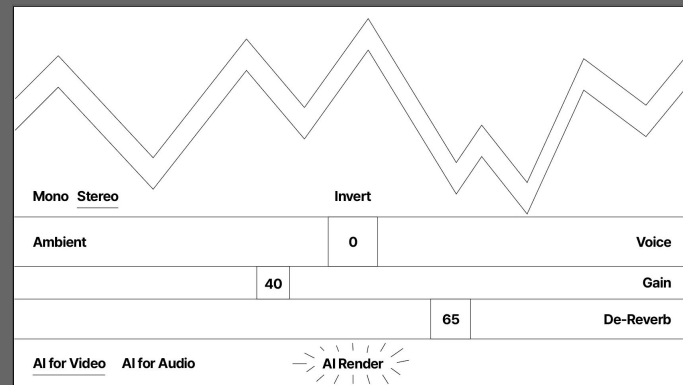
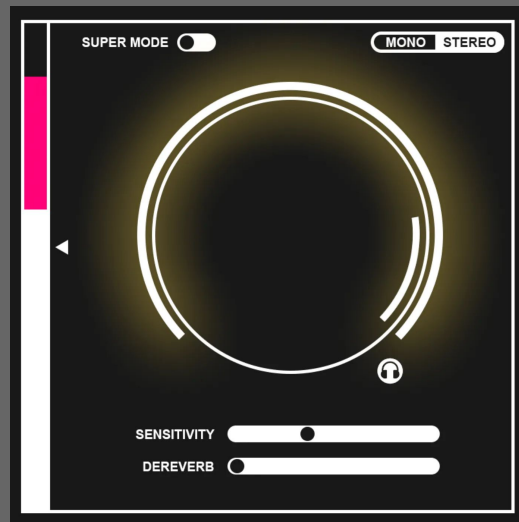
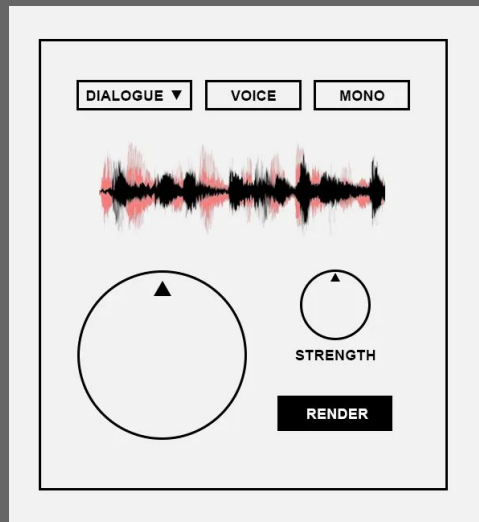
Writing a code
Writing a tests
Automating build pipeline
Supports variety of
platforms

Not Fun at all

Worrying about crack
Worrying about crash
Worrying about...
.
.
.

Developing a Plug-in

Definitely fun part



Developing a Plug-in

Definitely fun part

GOYO(beta) / Sep 2022



- Apply fully Physically Based 3D Rendering with OpenGL
- Implements dynamic deferred lighting.
- Large memory, GPU consumption

Supertone Clear / Sep 2023

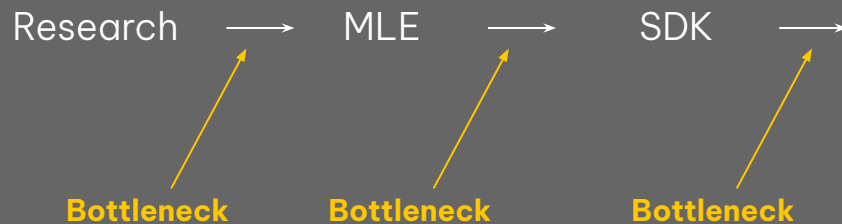


- Compact and restrained graphic.
- Focused on User Experience .
- More informative

Developing a Plug-in

Not Fun at all

Waiting for bottleneck...

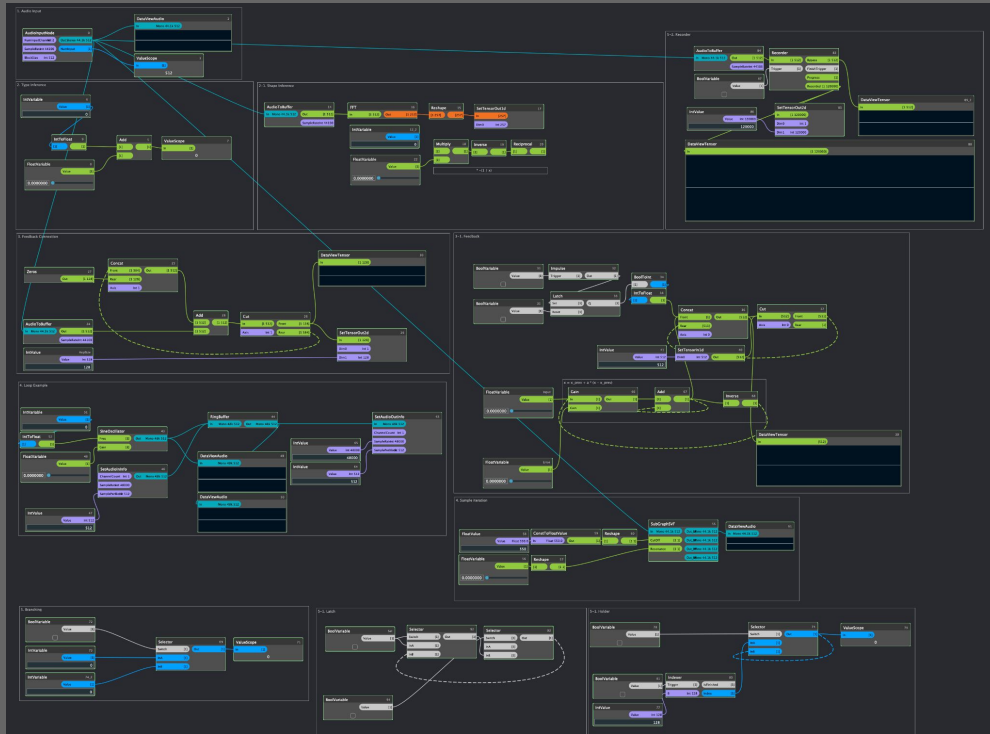




Deep Dive

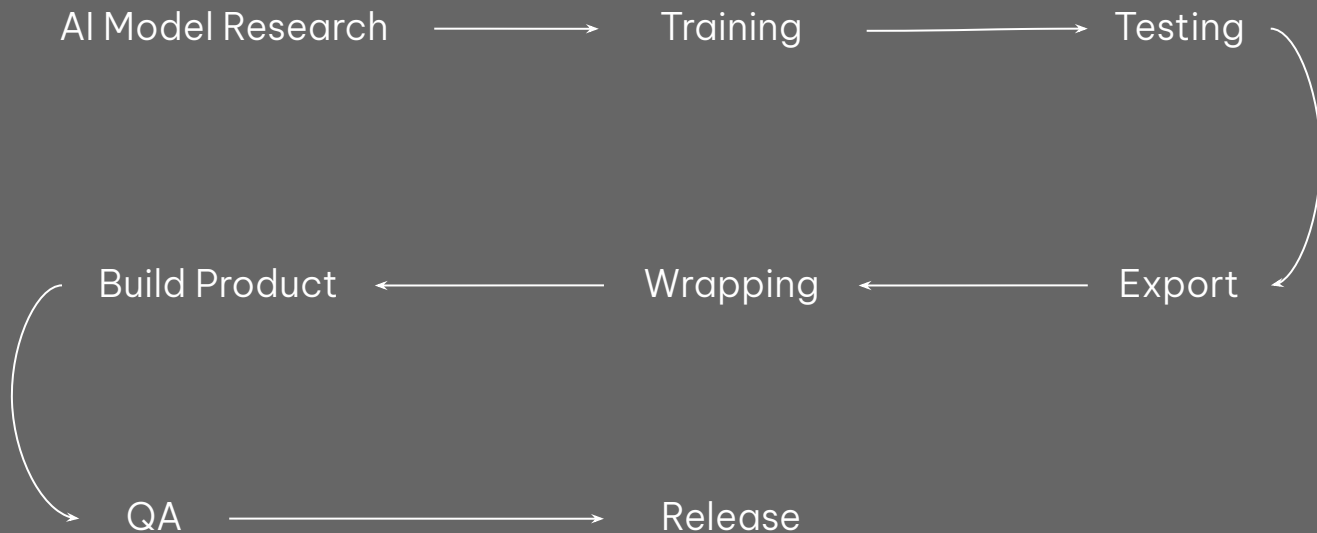
Naud Framework



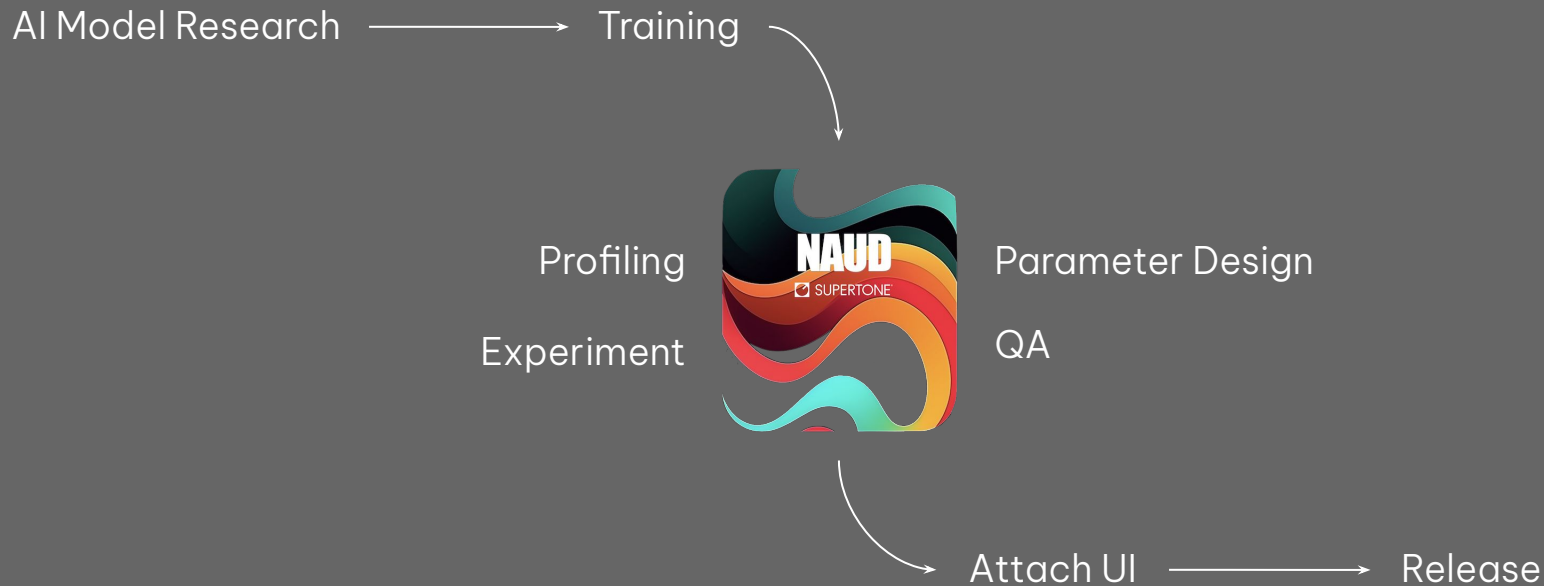


- Visual Programming Tool
- Inference Engine Integration
- Support Audio / MIDI IO
- Performance-Oriented
- Modular design

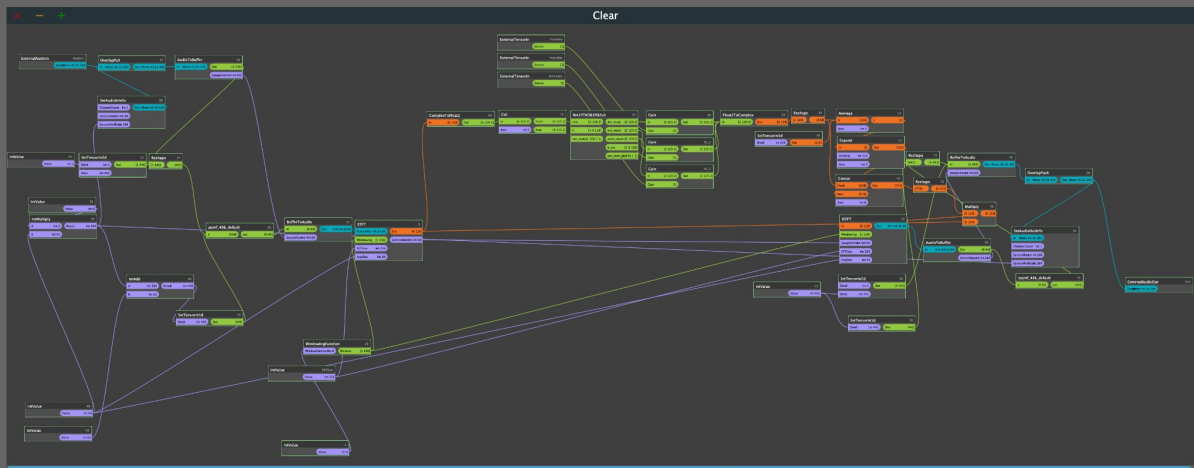
Launching an AI-powered audio software



How does it change the pipeline?



How Clear was made



| | |
|------------------|--------------------|
| Clear | 11 |
| AudMbm 44.1k 384 | Out Mono 44.1k 384 |
| VoiceDry | [1] |
| VoiceWet | [1] |
| Ambience | [1] |



Community Beta Release

Mar. 2024

Stay tuned.



Community Beta Release

Mar. 2024

Failed!

Stay tuned.



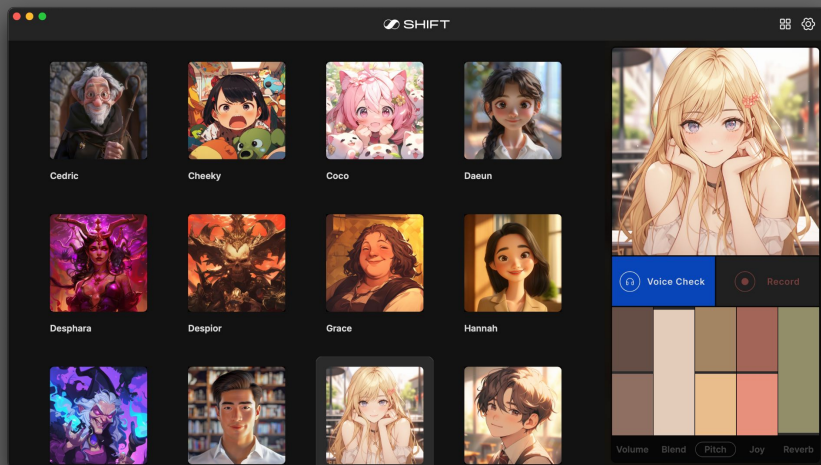
Deep Dive

Supertone Shift



Supertone Shift

Real-time voice conversion Software



- Desktop Application runs on Windows and macOS
- Can connect with Google Meet, Discord, Zoom etc.
- Provides variety of persona.
- Utilize DDSP to process real-time audio.
- Provides Pitch, Pitch variance and Blend.
- Runs on CPU in with 35ms latency.

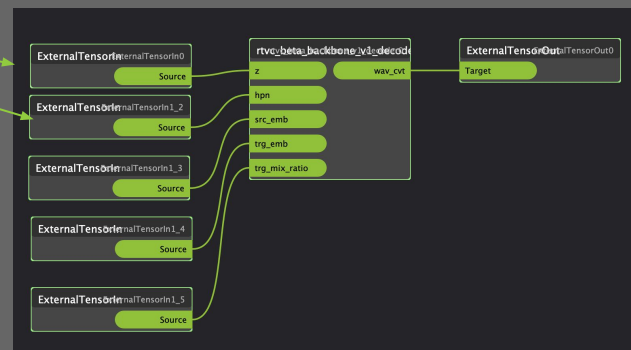
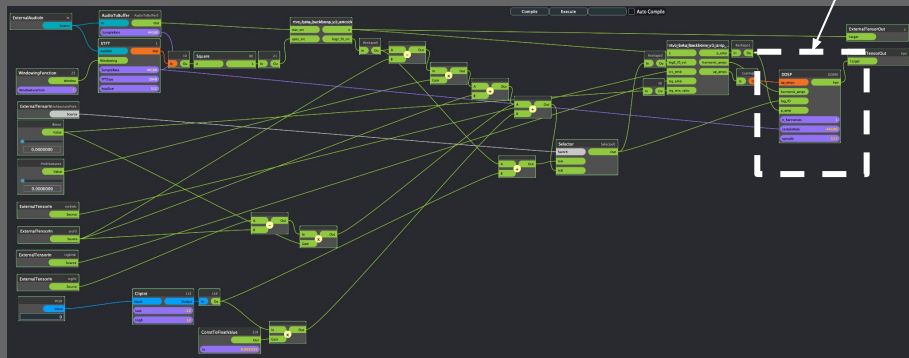
Developing Supertone Shift

With using Naud framework

Encoder

DDSP node

Decoder



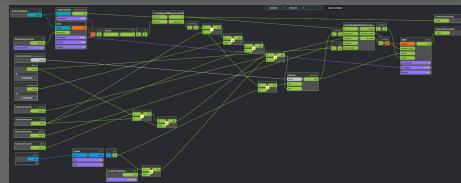
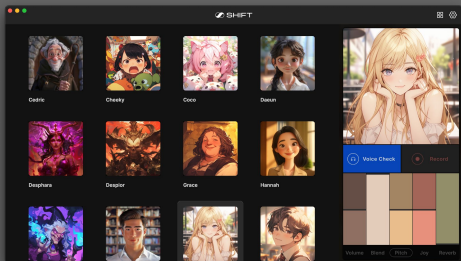
Input voice analysis

Output voice synthesis

Developing Supertone Shift

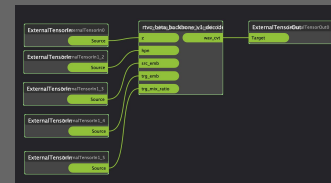
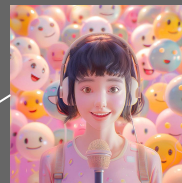
With using Naud framework

Encoder

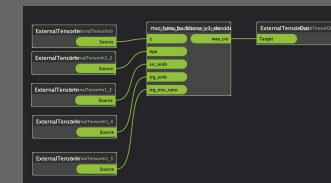
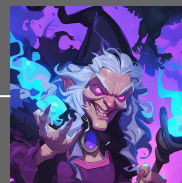


Input voice analysis

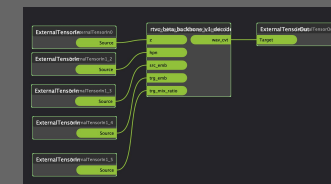
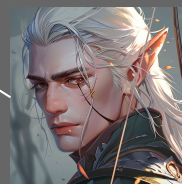
Decoder



Decoder



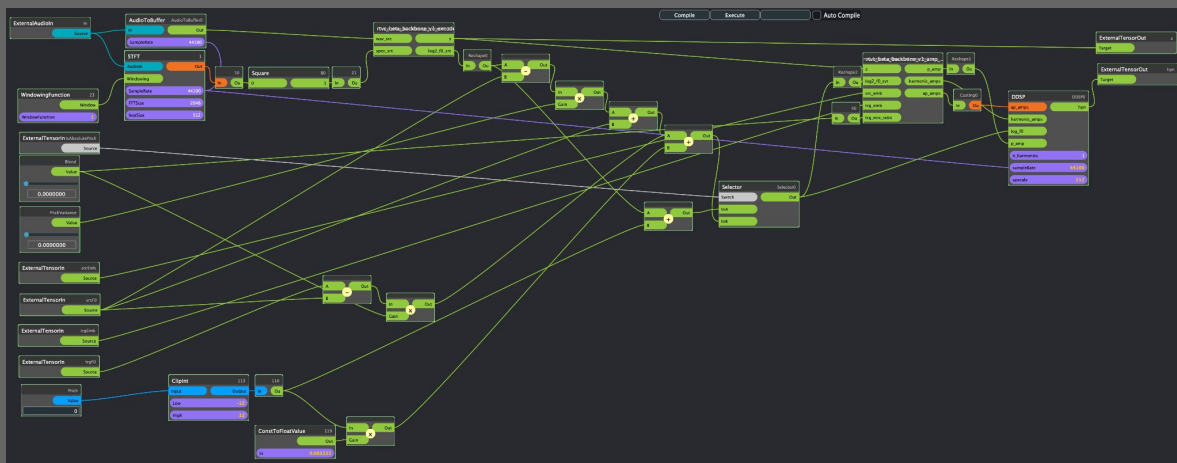
Decoder



Output voice synthesis

Developing Supertone Shift

With using Naud framework



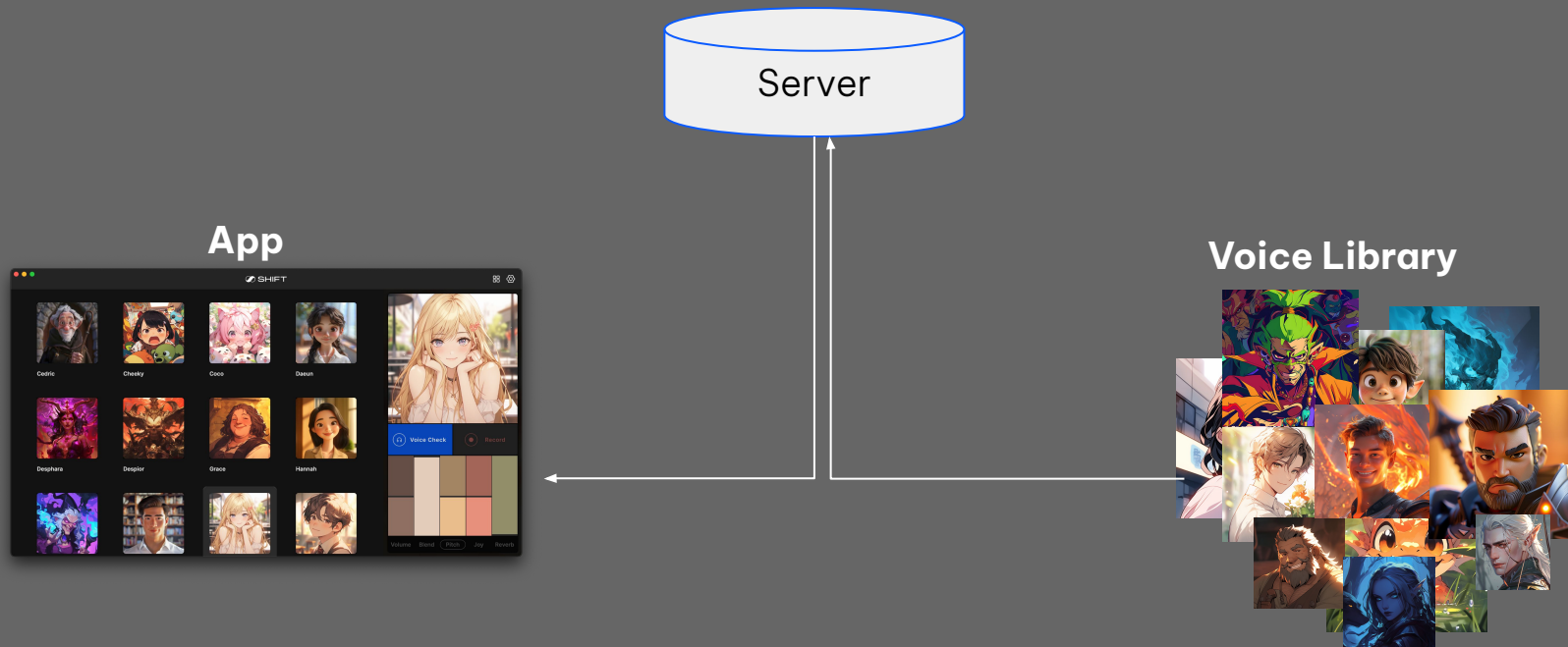
```

EncoderBlock
EncoderBlock | No Selection
1 <?xml version="1.0" encoding="UTF-8"?>
2
3 <Root>
4 <GraphInfo>
5 <AudioNode>
6 <STFT id="1"/>
7 <Abs id="10"/>
8 <Square id="80"/>
9 <WindowingFunction id="23" WindowFunction="1"/>
10 <Reshape id="21"/>
11 <Reshape id="40"/>
12 <ClipInt id="113" Low="-12" High="12"/>
13 <Casting id="116"/>
14 <Gain id="117"/>
15 <ConstToFloatValue id="119" In="0.083333"/>
16 <IntVariable id="Pitch"/>
17 <AudioToBuffer id="AudioToBuffer"/>
18 <ExternalTensorIn id="srcEmb"/>
19 <FloatVariable id="Blend"/>
20 <ExternalAudioIn id="In"/>
21 <ExternalTensorIn id="srcF0"/>
22 <Gain id="Gain"/>
23 <Subtract id="Subtract0"/>
24 <rtvc_beta_backbone_v1_amp_estimator
   id="rtvc_beta_backbone_v1_amp_estimator0"/>
25 <DPS0 id="DPS0" sampleRate="44100" upscale="512"/>
26 <rtvc_beta_backbone_v1_encoder id="rtvc_beta_backbone_v1_encoder0"/>
27 <ExternalTensorIn id="trgEmb"/>
28 <ExternalTensorIn id="trgF0"/>
29 <Reshape id="Reshape0"/>
30 <Reshape id="Reshape1"/>
31 <Casting id="Casting0"/>
32 <Reshape id="Reshape2"/>
33 <ExternalTensorOut id="Out"/>

```

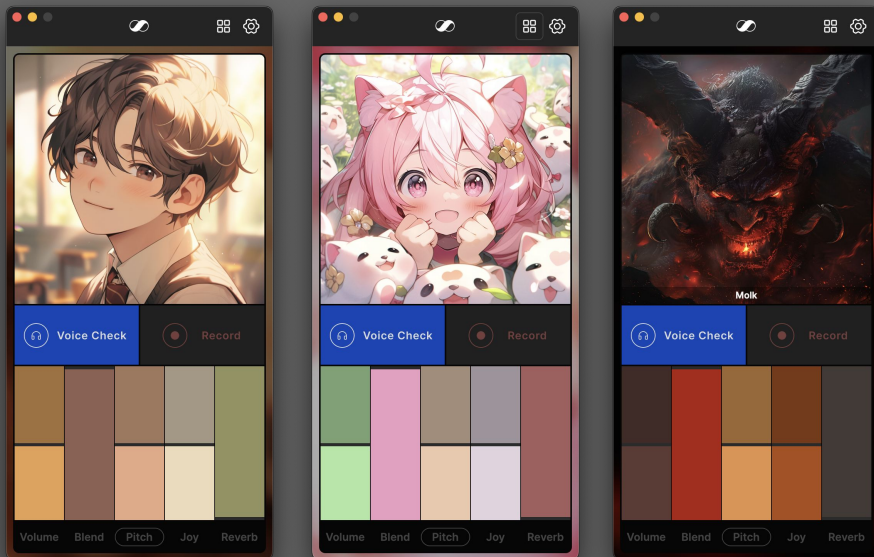
Developing Supertone Shift

With using Naud framework



Developing Supertone Shift

Who is user?



- Mostly non professional for audio.
 - VTuber
 - Podcaster
 - Voice contents creator
- Not familiar with traditional audio software.

Developing Supertone Shift

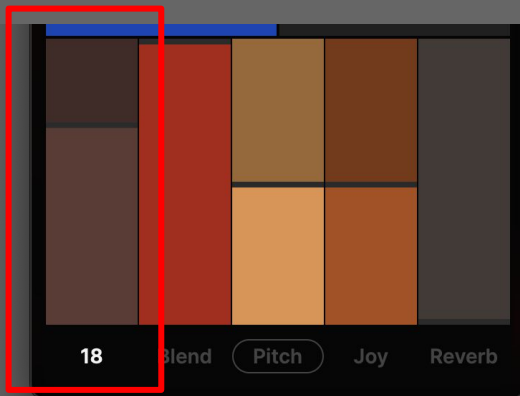
Volume slider problem

-24dB ~ +24dB

vs

1 ~ 25

- It is standard representation of volume.
- It is still intuitive representation.
- It is not that important thing, why should we argue with this?



- Why should our users have to know about standard?
- There is more intuitive representation.
- Because it matters whole experiences!



Deep Dive

Supertone Air



Supertone Air

Reverb & EQ Dialogue Match



Developing Supertone Air Implementing Transformer on Naud Framework

Yet Another Generative Model For Room Impulse Response Estimation *S, Lee et al*

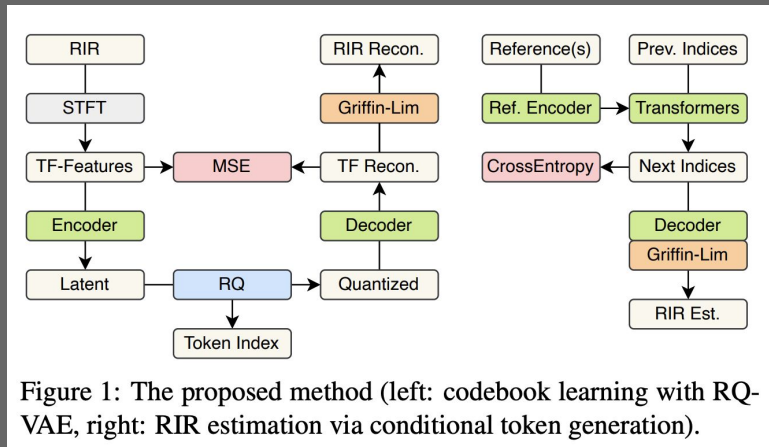
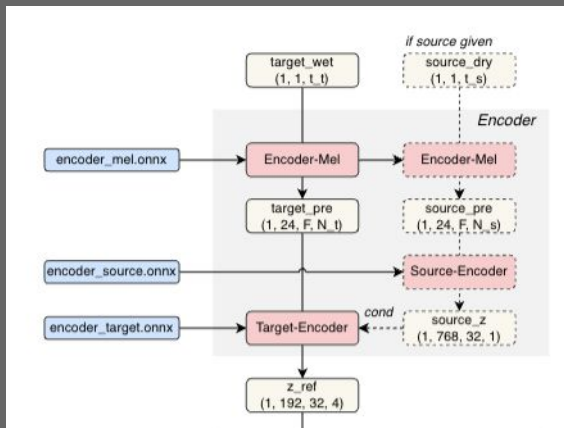


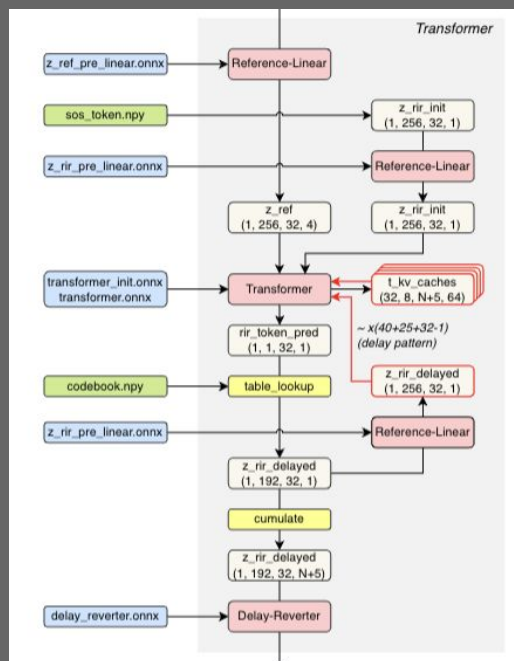
Figure 1: The proposed method (left: codebook learning with RQ-VAE, right: RIR estimation via conditional token generation).

Developing Supertone Air

Inference pipeline of AST model



Encoder



Transformer

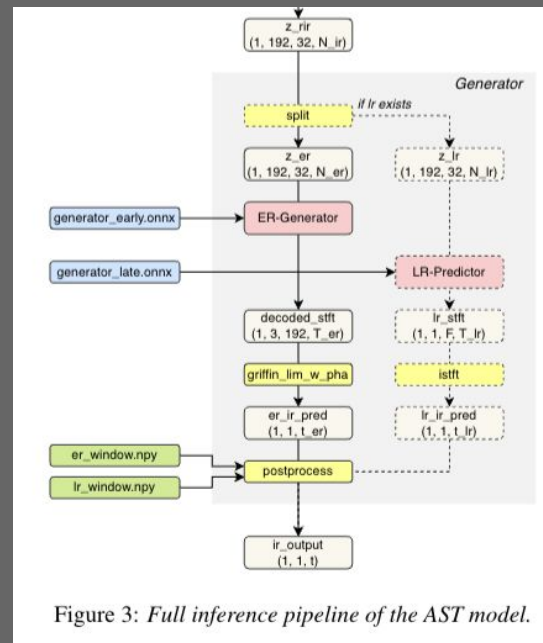
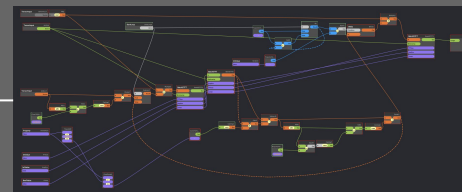
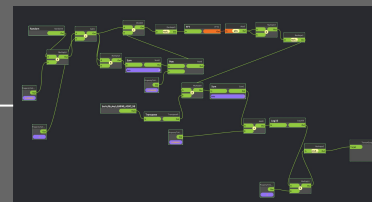
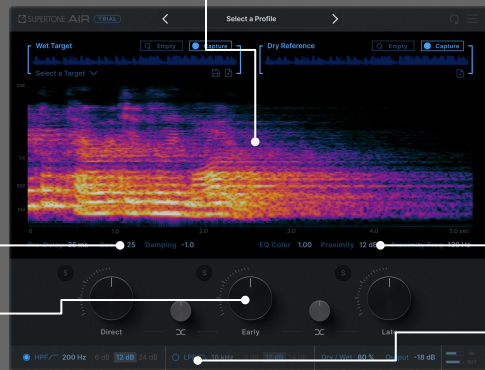
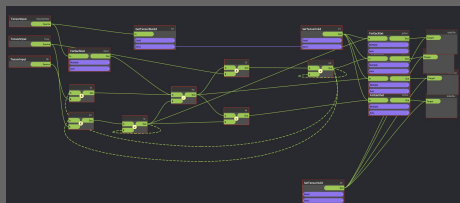
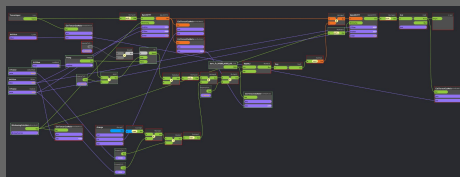


Figure 3: Full inference pipeline of the AST model.

Post-processor

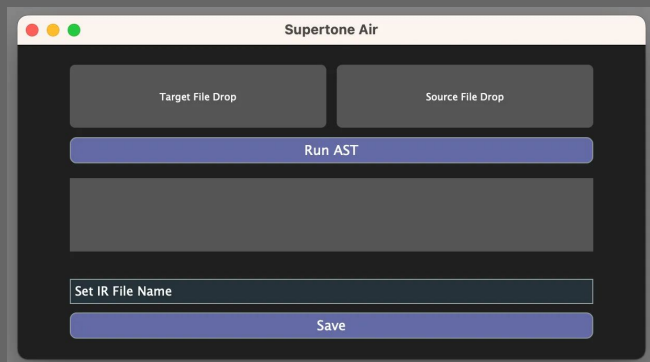
Developing Supertone Air

Implementing DSP algorithms with Naud Framework



Developing Supertone Air

Development timeline



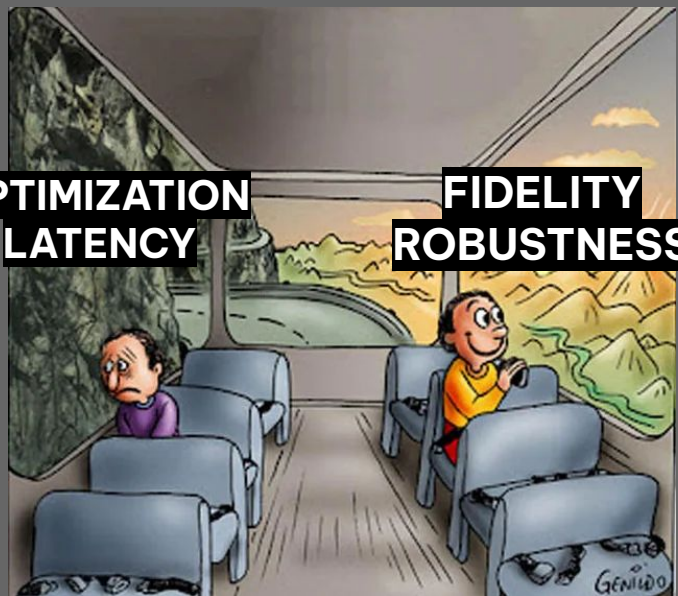
- Took less than 6 month including Market research, Product Design, Development.
- During the development, we took an advice from actual sound engineer to enhance the performance.

Developing Supertone Air

Different perspective

Performance!

- The new trained model has poorer performance.(Got slower)
- Let's just place the denoiser block before model input. It just makes simple.(To fit the deadline.)

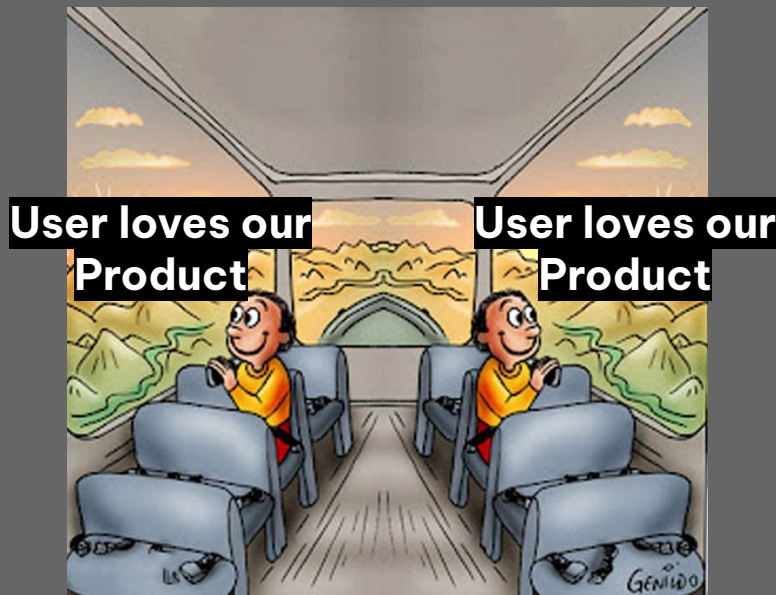


- It can't be! I've trained with more qualified data!(Got better fidelity)
- No. If the model has been well trained, than should handle input noise well. (which is not yet.)

Developing Supertone Air

Different perspective but same direction

Success!





Conclusions

Conclusions

- Build a good team.
 - Synergise you and your team's expertise.
 - Effective communication is always gold.
- Focusing on User, not Technology itself.
 - Look at the problems to be solved, the needs to be met.
 - Solution depends on who our customers are.
- Have Fun!

Thanks

Does anyone have any questions?

Email

rickysung@supertone.ai

Website

supertone.ai