# LIP-SYNC ML:

MACHINE LEARNING-BASED FRAMEWORK TO GENERATE LIP-SYNC ANIMATIONS IN FINAL FANTASY VII REBIRTH

© SQUARE ENIX
CHARACTER DESIGN: TETSUYA NOMURA / ROBERTO FERRARI

The second game of FINAL FANTASY VII remake trilogy

Previous title: FINAL FANTASY VII REMAKE (2020)

© SQUARE ENIX
CHARACTER DESIGN: TETSUYA NOMURA / ROBERTO FERRARI

- **Quality**
  - Express attractive characters and their dialogues.
  - Provide a game experience
    with the quality of our past CGI movie.
- **Efficiency**
  - Create a large amount of lip-sync animations
    only from audio.

SIGGRAPH 2024
DENVER+ 28 JUL – 1 AUG

FINAL FANTASY VII REMAKE

FINAL FANTASY VII REBIRTH

Phoneme-based method

→ Replace

ML-based method: Lip-Sync ML

Use machine learning approach.

- Input: Audio

- Output: Lip-sync animations

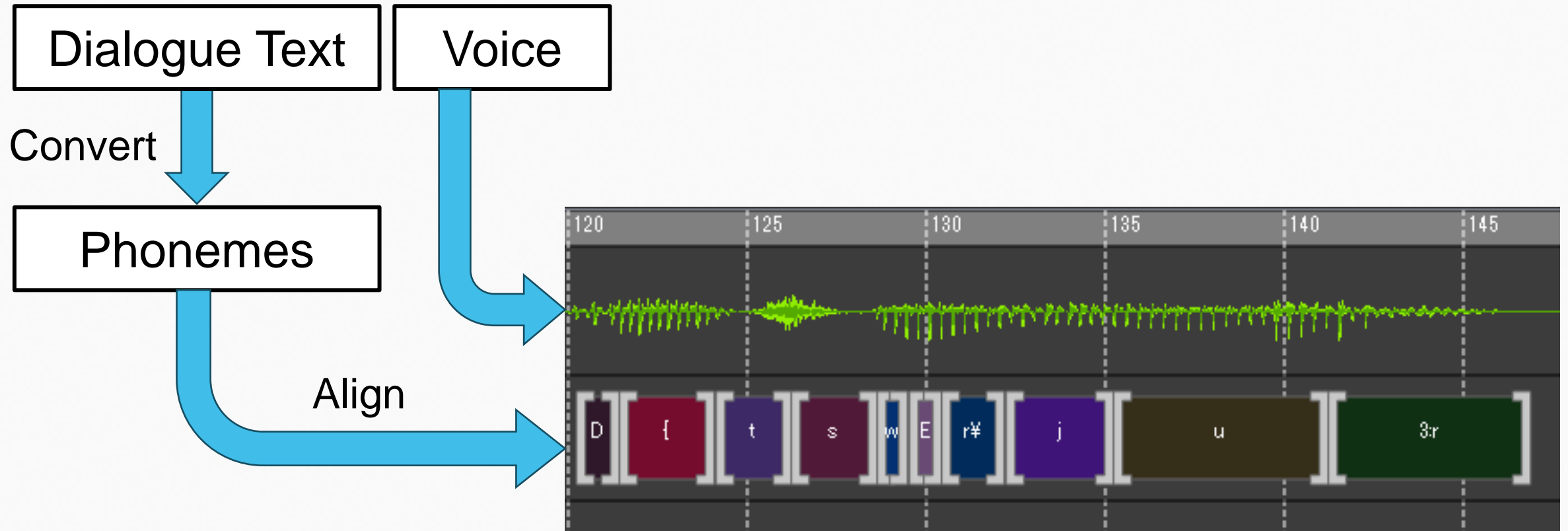- Training data: Cutscenes in FINAL FANTASY VII REMAKE

Dialogue Text

Voice

Convert

Phonemes

Align

Blend poses in the Lipmap corresponding to each phoneme.

7

default

aa

ee

fv

Lipmap is an asset
storing multiple mouth poses
for certain phonemes.

We choose 8 poses (incl. default)
for FINAL FANTASY VII REMAKE
and FINAL FANTASY VII REBIRTH.

- Input: Audio

- Output: Lip-sync animations

|  | **Bone-based animation** | **Pose weight animation (Lsdml)** |
|---|---|---|
| Feature | Better quality | Easy to synthesized with animations of other facial behaviors |
| Usage | Base keyframes to create lip-sync animations for cutscenes | Used in simple event scenes, combat and field actions |

Perform machine learning generation process.

Support batch processing to deal with multiple voices.

# CONVERT BONE ANIMATION INTO RIG CONTROLS

Determine rig parameters in each animation frame based on affected bones.

Edit Lsdml as animation curves in Maya if necessary.

In game engine, LsdmlPlayer plays the animation.

13

Lip-Sync ML
(Lsdml)

Previous
phoneme-based method

This table is about animations with voices

for simple event scenes and action sequences.

|  | **FINAL FANTASY VII REBIRTH** | **FINAL FANTASY VII REMAKE** |
|---|---|---|
| Method | Lip-Sync ML (Lsdml) | Phoneme-based method |
| The number of audio clips | About 136,000 | About 96,000 |
| The number of clips corrected manually | **10 ~ 20** | **More than 100** |

Make setting up the tool easier.

Users have to install necessary environment for machine learning currently.

SIGGRAPH 2024
DENVER+ 28 JUL — 1 AUG

LIP-SYNC ML: MACHINE LEARNING-BASED
FRAMEWORK TO GENERATE LIP-SYNC
ANIMATIONS IN FINAL FANTASY VII REBIRTH

# MACHINE LEARNING DETAILS

- About 3.5 hours of cutscene data from FINAL FANTASY VII REMAKE.

  o 53 different characters.

  o 3 facial skeletons: main characters, mob characters, Red XIII.

  o 2 languages: Japanese and English.

- Bone transform animation data (was easier to collect).

- Data augmentation: random speed and pitch changes.

18

- Designed as 2 sub-models, trained end-to-end.

- Voice model: transform audio into voice features.

- Animation model: generate animation data.

SIGGRAPH 2024
DENVER+ 28 JUL – 1 AUG

- Multiple independent styles: language, actor, rig.

- Each style has a set of exclusive values (Actor: Cloud, Aerith, Tifa...).

- Styles can also be left empty.

| Language | Japanese | English | |
|----------|----------|---------|------|
| **Actor** | Cloud | Aerith | Tifa |
| **Rig** | Main Chara | Mob Chara | RedXIII |

Language=None

Actor=Cloud

Rig=MainChara

- Designed as 2 sub-models, trained end-to-end.

- Voice model: transform audio into voice features.

- Animation model: generate animation data.

- Convert audio to mono.

- Resample to 19200 Hz (makes converting to 30 and 60 fps easier).

- Sample rate: 19200 Hz, window size: 200, stride: 160.

- Produces 120 spectrogram frames per second.

- Frequency range: 80 to 8000 Hz.

- Convert vertical axis to mel scale (logarithmic) using 256 bins.

- Closer to human perception of pitch.

- Changes in pitch become closer to vertical translations.

- Apply logarithm to the values (amplitude).

- Closer to human perception of volume.

- Can be seen as an image of 120 frames per second (time) x 256 mel bins (pitch), with a single value per pixel (volume).

- **Absolute pitch**: use each column directly as a feature vector (bag of features).

- **Relative pitch**: process as a 2D image (pitch invariance).

26

- Transfer information from height into feature depth.

- Progressively apply 2D convolutional networks with a stride of 2 in the vertical axis to halve height, while increasing feature depth.



| 2D Conv<br>depth=32<br>kernel=7x7<br>stride=1x2<br><br>ReLU<br>GroupNorm<br>groups=4 | 2D Conv<br>depth=48<br>kernel=7x3<br>stride=1x2<br><br>ReLU<br>GroupNorm<br>groups=4 | 2D Conv<br>depths=64<br>kernel=7x3<br>stride=1x2<br><br>ReLU<br>GroupNorm<br>groups=8 | 2D Conv<br>depth=96<br>kernel=7x3<br>stride=1x2<br><br>ReLU<br>GroupNorm<br>groups=8 | 2D Conv<br>depth=128<br>kernel=7x3<br>stride=1x2<br><br>ReLU<br>GroupNorm<br>groups=16 | Custom residual blocks |

- Extension of ResNet block [1].
- **GN**: Group Normalization [2].
- **SE**: Squeeze and Excitation [3].
- **Dotted**: ResNet shortcut connection.

28

- 3 groups of 3 consecutive residual blocks with different time dilations.

- In each group: 1st block halves height, 2nd and 3rd apply time dilation.

- Absolute pitch features are combined using one more residual block.

| 2D Residual Block (x 3) |
|---|
| depth=128 |
| kernel=7x3 |
| dilation=1,2,4 |
| stride=2,1,1 |

| 2D Residual Block (x 3) |
|---|
| depth=192 |
| kernel=7x3 |
| dilation=1,2,4 |
| stride=2,1,1 |

| 2D Residual Block (x 3) |
|---|
| depth=256 |
| kernel=7x3 |
| dilation=1,2,4 |
| stride=2,1,1 |

Reshape to 1D

+

1D Residual Block (x1)

depth=256
kernel=1

FC

Log mel spectrogram

Audio features

29

- 3x Transformer Encoder blocks [4].
  - Attend 1 second before / after.
  - Relative positional embeddings [5].
- Learn styles using LoRAs [6].
  - LoRAs for each style value.
  - Rank 8, added to qkv matrix.
  - Shared by all transformer blocks.

Styles

Audio features

LoRAs per style value

Transformer Encoder x 3

feature dims = 512
attention heads = 8
dropout = 0.1
mask = band matrix
1 second before/after

Relative positional embeddings

Voice features

- Voice features use 120 fps.

- Group and average voice features to 30 or 60 fps as needed.

- Apply one more transformer encoder block as before.

```
┌──────────────────┐      ┌──────────────────┐      ┌──────────────────┐      ┌──────────────────┐
│  Voice features  │─────▶│   Group and      │─────▶│                  │─────▶│  Frame features  │
│    (120 fps)     │      │    average       │      │  Transformer     │      │   (target fps)   │
└──────────────────┘      └──────────────────┘      │  encoder x 1     │      └──────────────────┘
┌──────────────────┐                                │                  │
│     Styles       │───────────────────────────────▶│                  │
└──────────────────┘                                └──────────────────┘
```

- Transform frame features into per-bone animation transform diffs.

- Format is the same as in training data (could also be rig parameters).

- Separate Fully Connected layers for each output rig / facial skeleton.

- Add character rest pose. Per frame outputs, no keyframes.

- L1 loss (difference abs).

Rest pose

L1 loss

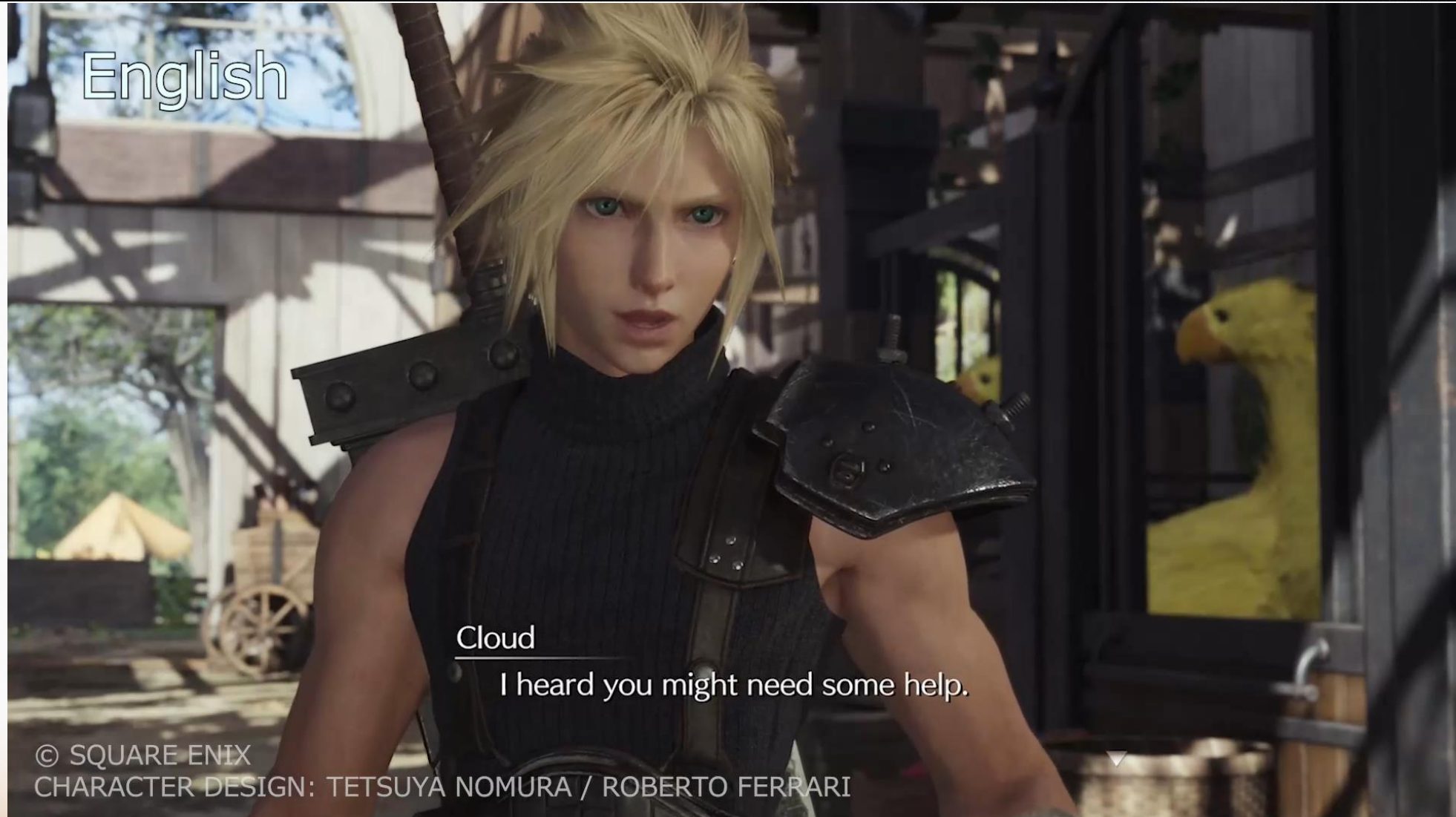| Frame features | FC for target rig | Animation diffs | Add pose | → | Animation (generated) | Animation (training data) |

- Transform frame features into lipmap pose weights (rig independent).

- Lipmap poses are fixed, pose weights are generated.

- Blending poses is a differentiable operation, can backpropagate.

- Same loss as animation outputs.

- Compensate how numerically relevant errors are for each value.

  - For example, a 0.5 error in a translation is likely small, but not in a scale factor.

  - Precompute the min-max ranges of each value in training data.

  - Rescale errors to make them relative to the sizes of their min-max ranges.

- Improve closing the mouth in generated animations.

  - Exploit the fact that rest poses have their mouths closed.

  - Scale up errors the closer training data values are to their rest pose.

In FINAL FANTASY VII REBIRTH,

our machine learning-based framework Lip-Sync ML

allowed animators to efficiently create high-quality lip-sync animations.

**This enabled characters to feel more natural when speaking, making dialogue scenes feel more immersive.**

# Thank you for your attention.

### Masato Nakada    Leandro Graciá Gil

nakmasat@square-enix.com    graclean@square-enix.com

Maya is a trademark or registered trademark of Autodesk, Inc.

All other trademarks are the property of the respective owners.

## For slides

(http://www.jp.square-enix.com

/tech/publications.html)

- [1] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. 2016.
  **Identity Mappings in Deep Residual Networks**. https://arxiv.org/pdf/1603.05027

- [2] Yuxin Wu and Kaiming He. 2018.
  **Group Normalization**. https://arxiv.org/pdf/1803.08494

- [3] Jie Hu, Li Shen, Samuel Albanie, Gang Sun, and Enhua Wu. 2019.
  **Squeeze-And-Excitation Networks**. https://arxiv.org/pdf/1709.01507

- [4] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones,
  Aidan N. Gomez, Łukasz Kaiser, and Illia Polosukhin. 2023.
  **Attention Is All You Need**. https://arxiv.org/pdf/1706.03762

- [5] Cheng-Zhi Anna Huang, Ashish Vaswani, Jakob Uszkoreit, Noam Shazeer, Ian Simon, Curtis Hawthorne,
  Andrew M. Dai, Matthew D. Hoffman, Monica Dinculescu, and Douglas Eck. 2018.
  **Music Transformer: Generating Music With Long-Term Structure**. https://arxiv.org/pdf/1809.04281

- [6] Edward Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, and Weizhu Chen. 2021.
  **LoRA: Low-Rank Adaptation of Large Language Models**. https://arxiv.org/pdf/2106.09685