

HumanDiffusion: diffusion model using perceptual gradients

Yota Ueda, Shinnosuke Takamichi, Yuki Saito, Norihiro Takamune, and Hiroshi Saruwatari
The University of Tokyo, Japan

1. SYNOPSIS

1. Proposes a diffusion model that represents a human-acceptable distribution.
2. Demonstrates that the proposed algorithm successfully samples data that follow a human-acceptable distribution.

2. RESEARCH BACKGROUND

Human-acceptable distribution

- A probability distr. in which data humans can accept as natural.
- A human-acceptable distribution is typically wider than a real-data distribution.



Modeling the distribution with DNN (deep neural network)

- Can generate speech in a human-acceptable distribution.
- Can generate speech that do not exist.

3. RELATED WORKS

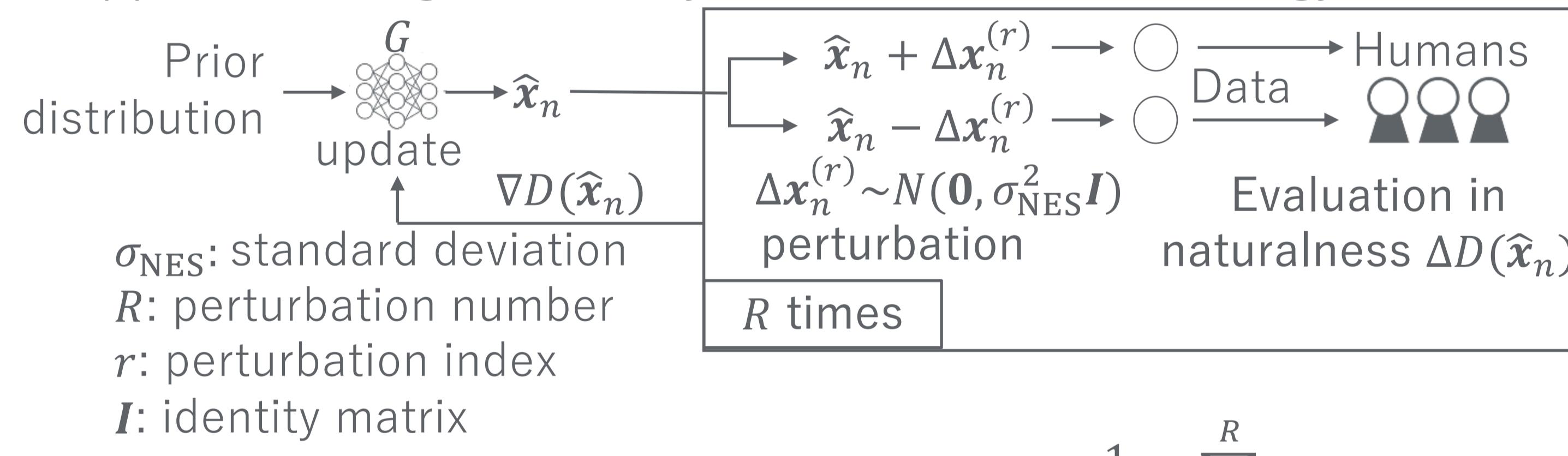
GAN (generative adversarial network)^[1]

- Consists of a DNN-generator and a DNN-discriminator.
- The generator outputs data in a **real-data distribution**.
- Gradients ∇D for updating generator's parameters **can** be calculated analytically.



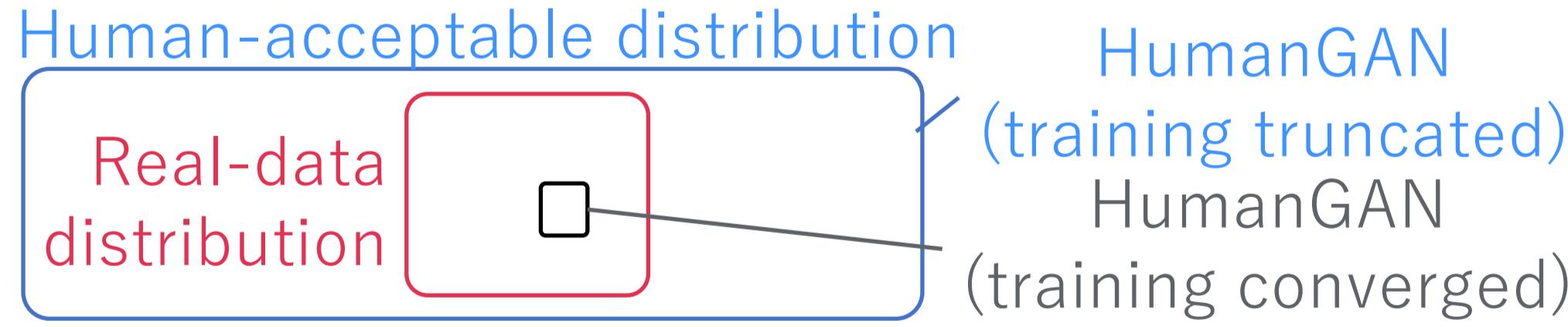
HumanGAN^[2]

- Consists of a DNN-generator and a humans-discriminator.
- The generator outputs data in a **human-acceptable distribution**.
- Gradients ∇D **cannot** be calculated analytically.
- Approximates gradients by natural evolution strategy (NES)^[4].



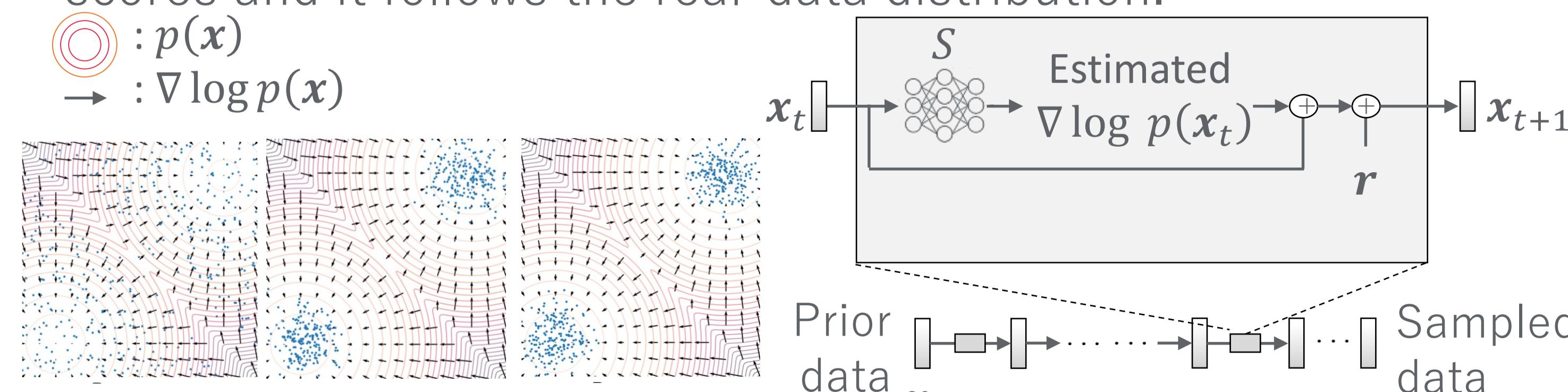
- Gradients are approximated as $\nabla D(\hat{x}_n) \approx \frac{1}{2\sigma_{NES}R} \sum_{r=1}^R \Delta D(\hat{x}_n^{(r)}) \cdot \Delta x_n^{(r)}$

- Problem: needs heuristic training truncation



Diffusion model^[3]

- Consists of a score network S only.
- The score network models score $\nabla \log p(\mathbf{x})$ of the **real-data distribution $p(\mathbf{x})$** .
- Data is sampled with Langevin dynamics using estimated scores and it follows the real-data distribution.



- From: <https://yang-song.net/blog/2021/score/>
- Sampling with Langevin dynamics
- $$\mathbf{x}_{t+1} = \mathbf{x}_t + \frac{\epsilon^2}{2} S(\mathbf{x}_t) + \epsilon \mathbf{r}$$
- If score of humans' evaluation function $D(\mathbf{x})$ is modeled instead of $p(\mathbf{x})$, sampled data follows a human-acceptable distribution.

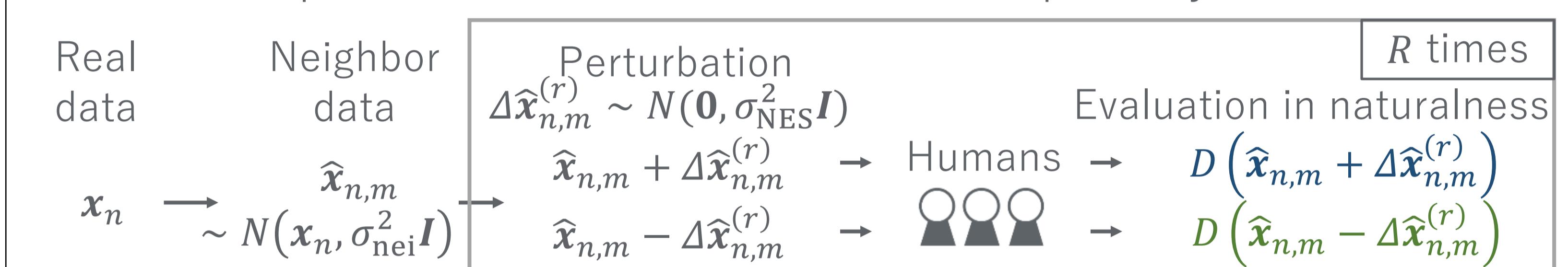
4. PROPOSED METHOD: HumanDiffusion

Method

- Sampling data $\mathbf{x} \sim \frac{1}{\int D} D(\mathbf{x})$ with score $\nabla \log D(\mathbf{x})$

Approximation of score

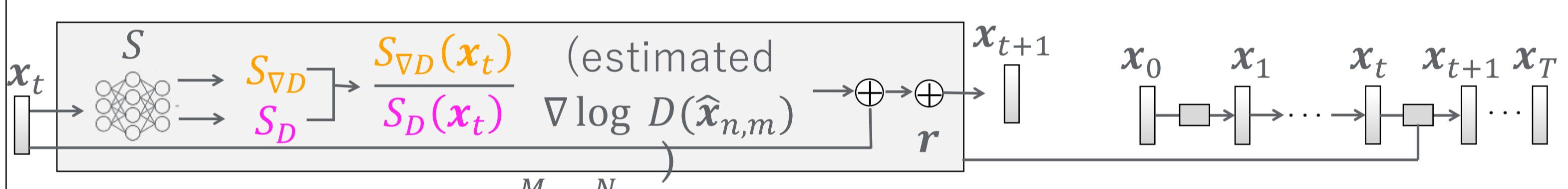
- Approximating score because it cannot be calculated analytically.
- Neighbor data is sampled from real data for gradient around them.
- Evaluates perturbed data in naturalness acceptability for NES.



- Score is calculated as

$$\nabla \log D(\hat{x}_{n,m}) = \frac{\nabla D(\hat{x}_{n,m})}{D(\hat{x}_{n,m})} \quad \nabla D(\hat{x}_{n,m}) = \frac{1}{2\sigma_{NES}R} \sum_{r=1}^R (D(\hat{x}_{n,m} + \Delta \hat{x}_{n,m}^{(r)}) - D(\hat{x}_{n,m} - \Delta \hat{x}_{n,m}^{(r)})) \cdot \Delta \hat{x}_{n,m}^{(r)}$$

Training score network and sampling



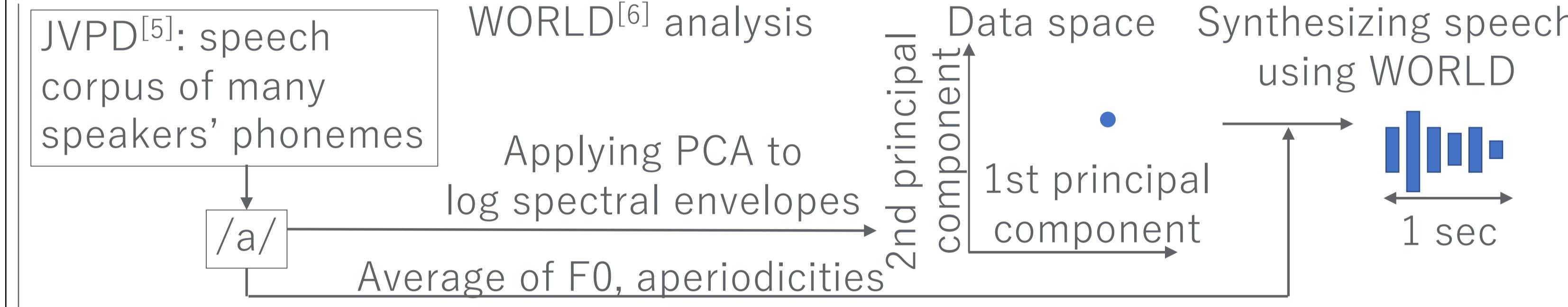
- Loss function: $L = \sum_{m=1}^M \sum_{n=1}^N (|S_{VD}(\hat{x}_{n,m}) - \nabla D(\hat{x}_{n,m})|^2 + |S_D(\hat{x}_{n,m}) - D(\hat{x}_{n,m})|^2)$
- Sampling: $x_{t+1} = x_t + \frac{\epsilon^2}{2} \frac{S_{VD}(x_t)}{S_D(x_t)} + \epsilon r$

5. EXPERIMENTAL EVALUATIONS

Purpose

- Confirm HumanDiffusion models human-acceptable distribution.

Preparation of data space

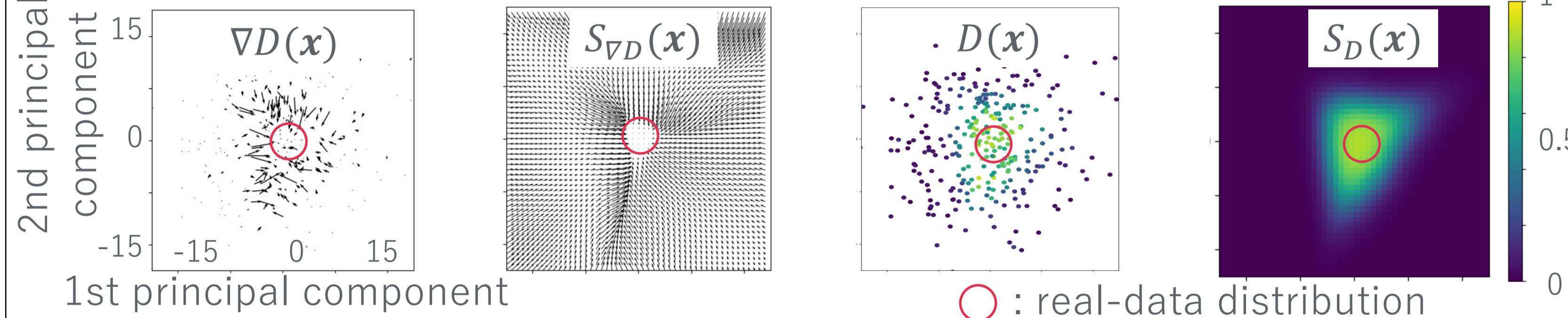


Hyper-parameters

- $n: 100, m: 3, \#humans: 105, \sigma_{nei}: 10, \sigma_{NES}: 1.0, \epsilon: 0.003, T: 10000, R: 20$
- Score network: 3 fully-connected layers, optimizer: Adam

Evaluation 1: confirm score network learns $D(\mathbf{x})$ and $\nabla D(\mathbf{x})$

- $S_D(\mathbf{x})$ got higher values around the real-data distribution.
- $S_{VD}(\mathbf{x})$ pointed to the real-data distribution.
- **Score network was trained successfully.**



Evaluation 2: perceptual evaluation of sampled data

- Humans evaluated sampled data and real data in naturalness.
- The evaluation value was almost same although the distribution of sampled data was wider than real-data distribution.
- **HumanDiffusion samples data in human-acceptable distribution.**

