# END-TO-END FEEDBACK LOSS IN SPEECH CHAIN FRAMEWORK VIA STRAIGHT-THROUGH ESTIMATOR

Andros Tjandra[1], Sakriani Sakti[1,2], Satoshi Nakamura[1,2]

1) Nara Institute of Science and Technology, Japan   2) RIKEN, Center for Advanced Intelligence Project AIP, Japan

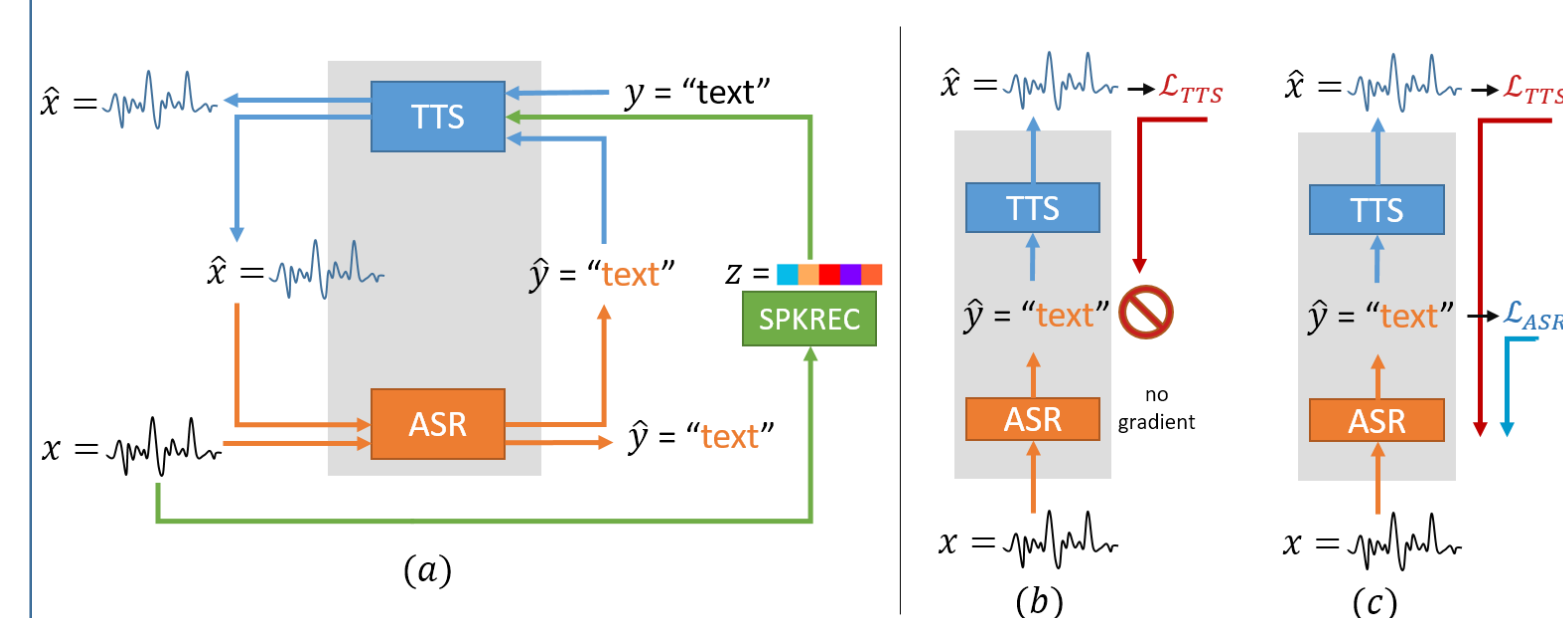{andros.tjandra.ai6, ssakti, s-nakamura}@is.naist.jp

## 1. Introduction

- Speech chain model integrates ASR and TTS into a single cycle during training.
- By combining both models, we could train with auxiliary feedback loss.
- Problem:
  - The output from ASR are discrete tokens
  - Non-differentiable (ASR → TTS)
- Solution:
  - Apply straight-through estimator on Gumbel-softmax or $argmax$ sample
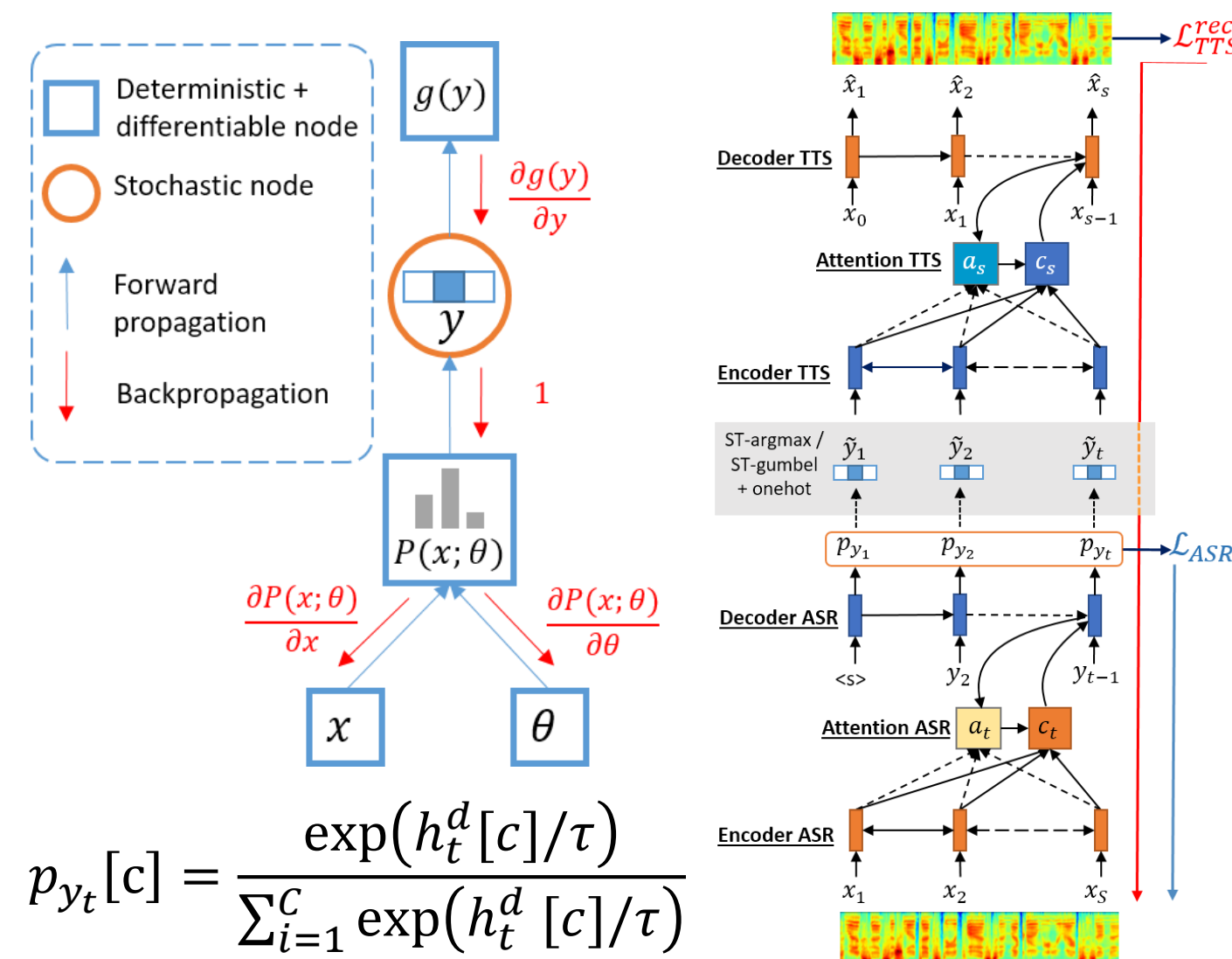
## 2. Speech Chain and Feedback loss



Feedback loss: $\mathcal{L}_{TTS} = || x - \hat{x}||_2^2$ where $x = TTS(\hat{y}, z)$

a) Speech chain loop with speaker embedding module.
b) Original: feedback $\mathcal{L}_{TTS}$ can't be backpropagated through variable $\hat{y}$.
c) **Proposal**: Estimate gradient through variable $\hat{y}$ with straight-through estimator.

## 3. Straight-through Estimator (ST)



$$p_{y_t}[c] = \frac{\exp(h_t^d[c]/\tau)}{\sum_{i=1}^{C} \exp(h_t^d[c]/\tau)}$$

$\tau$ = temperature

**a) ST-argmax**
Deterministic choosing token by highest probability.

$$\tilde{y}_t = argmax_c \, p_{y_t}[c]$$

**b) ST-Gumbel softmax**
Sampling a token from $p_{y_t}[c]$:

$$p_{y_t}[c] = \frac{\exp((h_t^d[c] + g_c)/\tau)}{\sum_{i=1}^{C} \exp((h_t^d[c] + g_c)/\tau)}$$

$$\tilde{y}_t \sim Cat(p_{y_t}[1], \ldots, p_{y_t}[C])$$

New gradient $\mathcal{L}_{TTS}$ w.r.t. $\theta_{ASR}$

$$\frac{\partial \mathcal{L}_{TTS}^{rec}}{\partial \theta_{ASR}} = \sum_{t=1}^{T} \frac{\partial \mathcal{L}_{TTS}^{rec}}{\partial \tilde{y}_t} \cdot \frac{\partial \tilde{y}_t}{\partial p_{y_t}} \cdot \frac{\partial p_{y_t}}{\partial \theta_{ASR}}$$

$$\approx \sum_{t=1}^{T} \frac{\partial \mathcal{L}_{TTS}^{rec}}{\partial \tilde{y}_t} \cdot \mathbb{1} \cdot \frac{\partial p_{y_t}}{\partial \theta_{ASR}}.$$

## 4. Experiment

- Features: log Mel-spec (80-dim)
- Text: 26 letters (A-Z)+(',- ) + <noise>
- Dataset: Wall Street Journal
  - Train: train_si284 (WSJ1)
  - Dev: dev93 & Test: eval92
- Hyperparams $\tau = [0.25, 0.5, 1, 2]$

Result on WSJ-1

| Baseline | | | |
|---|---|---|---|
| **Model** | | **CER (%)** | |
| Att MLP | | 7.12 | |
| Att MLP-MA | | 6.43 | |
| **Proposed** | | | |
| **Model** | **Generation** | **ST** | **CER (%)** |
| Att MLP-MA | Teacher-forcing | argmax | 5.75 |
| Att MLP-MA | Teacher-forcing | gumbel | 5.7 |
| Att MLP-MA | Greedy | argmax | 5.84 |
| Att MLP-MA | Greedy | gumbel | 5.88 |

## 5. Discussion

- We introduced ST-estimator for training ASR module based on TTS feedback loss.
- The gradient in discretization problem can be replaced by identity Jacobian matrix.
- Our experiment shows that by adding auxiliary feedback loss, we improve the ASR performance further.