

Tencent AI Lab

SANDGLASSET: A LIGHT MULTI-GRANULARITY SELF-ATTENTIVE NETWORK FOR TIME-DOMAIN SPEECH SEPARATION

Max W. Y. Lam*

Jun Wang*

Dan Su*

Dong Yu†

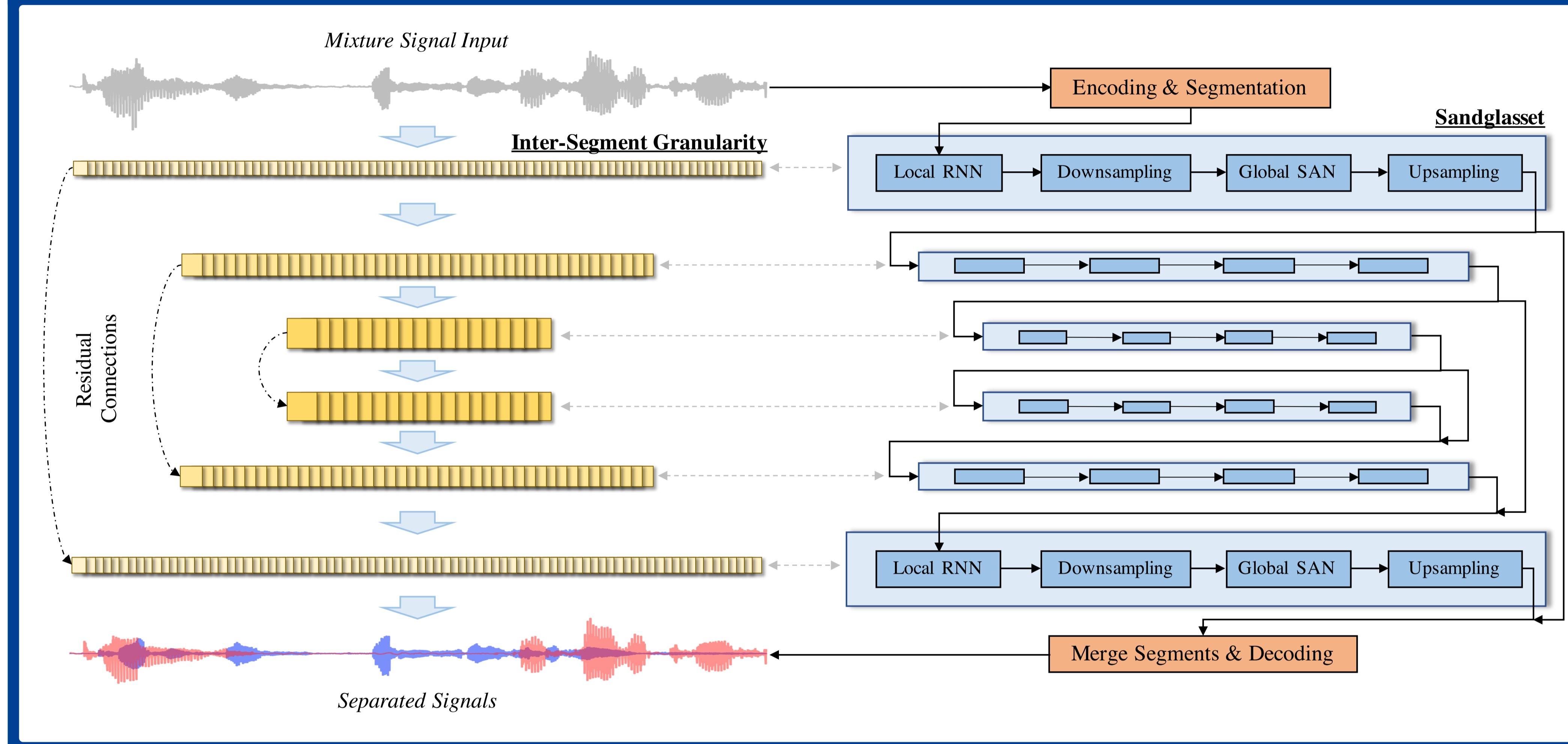
* Tencent AI Lab, Shenzhen, China

† Tencent AI Lab, Bellevue WA, USA

INTRODUCTION

- Speech separation is a fundamental component for many downstream speech processing tasks.
- Single-channel speech separation has recently been advanced by the time-domain audio separation networks (TasNets) (Luo and Mesgarani, 2018).
- State-of-the-art (SOTA) models both employ a dual-path technique, which is to process the segment sequence in an intra-segment (local) direction and an inter-segment (global) direction alternatively.
- In our previous work of GALR, we found that self-attentive networks (SANs) are superior over RNNs in modeling the inter-segment sequence.
- SAN can connect every element to another element with a direct path (i.e., in $\mathcal{O}(1)$ time), in contrast to $\mathcal{O}(N)$ time in RNNs.

SANDGLASSET



RESULTS

1) Performances on WSJ0-2mix:

Model	Params.	SI-SNRi	SDRi
BLSTM-TasNet	23.6M	13.2	13.6
Conv-TasNet	8.8M	15.3	15.6
DPRNN	2.6M	18.8	19.1
DPTNet	2.7M	20.2	20.6
Sandglasstet (w/o RES)	2.3M	20.1	20.3
Sandglasstet (SG)	2.3M	20.3	20.5
Sandglasstet (MG)	2.3M	20.8	21.0
Sandglasstet (MG) + PT	2.3M	21.0	21.2
Gated DPRNN + Spk ID	7.5M	20.1	-
Wavesplit + Spk ID	†42.5M	21.0	21.2

2) Performance on WSJ0-3mix:

Model	Params.	SI-SNRi	SDRi
Conv-TasNet	8.8M	12.7	13.1
DPRNN	2.6M	14.7	-
Sandglasstet (MG)	2.3M	17.1	17.4
Gated DPRNN + Spk ID	7.5M	16.7	-
Wavesplit + Spk ID	†42.5M	17.3	17.6

3) Cost Analysis:

Model	Params.	Memory (GB)	GFLOPs (10^9)
DPRNN	2.6M	1.97	84.7
Sandglasstet	2.3M	0.82 (↓58.4%)	28.8 (↓66.0%)

MULTI-GRANULARITY

- **Existing Method:** Use a fixed segment size unchanged throughout all layers; The global contexts being learnt in the inter-segment sequence is limited to one manually defined scale.
- **Fact:** Speech signals contain different level of contexts, for example, the phoneme scale, the syllable scale, and the word scale, what we refer to as the multi-granularity.
- **Our Observation:** SANs have superior capabilities in modeling sequences of high-level contexts, as examined in LM and in NLP.
- **Our Idea:** Design a novel network that allows SANs to capture multi-granularity information for enhancing contextual modeling and computational efficiency.
- **Our Proposed:** Sandglasstet, named for its sand-glass shape and its modest model size and complexity.

SANDGLASSET BLOCK

- **Sandglasstet block:** Each consists of these modules: 1) A **local RNN**; 2) A **global SAN**; 3) A **down-sampling** and an **upsampling** modules for changing the context granularity.
- For the b -th block, $\mathcal{X}_b \in \mathbb{R}^{D \times K \times S}$ is the block input enclosing S segments each containing K -length D -dimensional features.
- Mathematically, a Sandglasstet block computes the following:

$$\mathcal{Y}_b^{LR} = [\text{Linear}(\text{RNN}_b(\mathcal{X}_b[:, :, s])), s = 1, \dots, S], \quad (1)$$

$$\mathcal{Y}_b^{GA} = \text{US}_b(\text{SAN}_b(\text{DS}_b(\text{LN}(\mathcal{Y}_b^{LR}) + \mathcal{X}_b))), \quad (2)$$

$$\text{SAN}(\mathcal{X}) = [\text{SelfAttn}(\text{LN}(\mathcal{X}[:, k, :]) + \mathbf{P}), k = 1, \dots, K], \quad (3)$$

- The downsampling and upsampling operations are defined as

$$\text{DS}_b(\mathcal{X}) = \begin{cases} \text{Conv1D}_K(\mathcal{X}; 4^b) & \text{if } b \leq N/2; \\ \text{Conv1D}_K(\mathcal{X}; 4^{N-b-1}) & \text{if } b > N/2. \end{cases} \quad (4)$$

$$\text{US}_b(\mathcal{X}) = \begin{cases} \text{ConvTrans1D}_K(\mathcal{X}; 4^b) & \text{if } b \leq N/2; \\ \text{ConvTrans1D}_K(\mathcal{X}; 4^{N-b-1}) & \text{if } b > N/2, \end{cases} \quad (5)$$

where

$$\text{Conv1D}_K(\mathcal{X}; \tau) \in \mathbb{R}^{D \times [K/\tau] \times S} \quad (6)$$

$$\text{ConvTrans1D}_K(\mathcal{X}; \tau) \in \mathbb{R}^{D \times K\tau \times S} \quad (7)$$

CONCLUSIONS

- This paper proposes a novel network named Sandglasstet for time-domain speech separation.
- Sandglasstet applies a downsampling-upsampling mechanism to the global SAN for modeling multi-granularity contexts.
- As the smallest TasNet in size, Sandglasstet achieved the state-of-the-art results on two benchmark datasets and is also low-cost in terms of memory and computations.