HYPERCONNECT

1. Motivation

- To deploy keyword spotting (KWS) on mobile devices, it should be **fast** enough and **accurate** for real-time inference.
- However, previous studies are **either** fast or accurate.

Model	Acc. (%)	Time (<i>ms</i>)	FLOPs	Params
CNN-2	84.6*	1.2	1.5M	148K
Res15	95.8 *	424	1950.0M	239K

- Problem of 2D convolution.
 - Conventionally, MFCC is considered as input 2D image.
 - **Receptive fields** of conventional CNN are not enough to cover all features at specific time point.
 - It needs deeper or wider layers to cope with all frequency domains, which **slows down** the model.



2. Temporal Convolution



- Large receptive field of audio features.
- Enable a model to achieve better accuracy even with less computation.
- Small footprint and low computational complexity.
- Suitable for fast inference.

Temporal Convolution for Real-time Keyword Spotting on Mobile Devices

Seungwoo Choi^{*}, Seokjun Seo^{*}, Beomjun Shin^{*}, Hyeongmin Byun, Martin Kersner, Beomsu Kim, Dongyoung Kim⁺, Sungjoo Ha⁺



K



3. Proposed model



Block (s=1)

Block (s=2)

4. Experimental Results

Model	Acc. (%)	Time (<i>ms</i>)	FLOPs	Params
CNN-1	90.7	32	76.1M	524K
CNN-2 DS-CNN-S DS-CNN-M DS-CNN-L Res8-Narrow	84.6	1.2	1.5M	148K
	94.4	1.6	5.4M	24K
	94.9	5.2	19.8M	140K
	95.4	16.8	56.9M	420K
	90.1	47	143.2M	20K
Res8	94.1	174	795.3M	111K
Res15-Narrow	94.0	107	348.7M	43K
Res15	95.8	424	1950.0M	239K
TC-ResNet8	96.1	1.1	3.0M	66K
TC-ResNet8-1.5	96.2	2.8	6.6M	145K
TC-ResNet14	96.2	2.5	6.1M	137K
TC-ResNet14-1.5	96.6	5.7	13.4M	305K

* The accuracy is calculated on Google Speech Commands dataset and the speed is measured on Google Pixel 1.

TC-ResNet

5. Conclusion



Paper

- Accuracy improven
- Compared to the pr TC-ResNet8 show
- Compared to the pr TC-ResNet8 show
- TC-ResNet can adju
- By modulating th
- Ablation study show for fast and accurate
 - 2D-ResNet8 2D-ResNet8-Pool

 - layer to reduce computation.



• State-of-the-art performance on Google Speech Commands. • In terms of both speed and accuracy.

• By adopting **temporal convolution**.

• TC-ResNet facilitates real-time KWS on mobile devices.

• Not estimating by FLOPs but measuring on mobile devices.

• Implementation and benchmark tools are publicly released. https://github.com/hyperconnect/TC-ResNet



Github

nent regardless of the speed.							
revious most accurate model.							
vs 385x faster inference speed.							
revious fastest model.							
vs 11.5%p accuracy improvement.							
ist tradeoff between accuracy and speed.							
ne number of layers and width multiplier (k).							
ws temporal convolution is a key component							
ze KWS.							
96.1	1	0.1	35.8M	66K			
94.9)	3.5	4.0M	66K			

 2D-ResNet8: temporal convolution -> 2D convolution. 2D-ResNet8-Pool: + additional pooling after first convolution