# Few-Shot Learning Neural Network for Audio-Visual Speech Enhancement

## Maya Rapaport

M.Sc. Seminar
Supervisor: **Prof. Israel Cohen**

Department of Electrical Engineering, Technion, Israel Institute of Technology

March 2021

# Outline

Background & Motivation

**Audio-Visual Speech Enhancement**
Problem: Speaker Dependency
Proposed Algorithm

**Few-Shot Learning**
Problem: Dependency on the shots
Proposed Algorithm

Conclusions & Future Work

# Research Contributions

**1.** Overcoming speaker dependency for real-time mobile applications.

**Proposing**
**Fast Adaptation Speech Enhancement (FASE) model,**
**Inspired by few-shot learning methods**

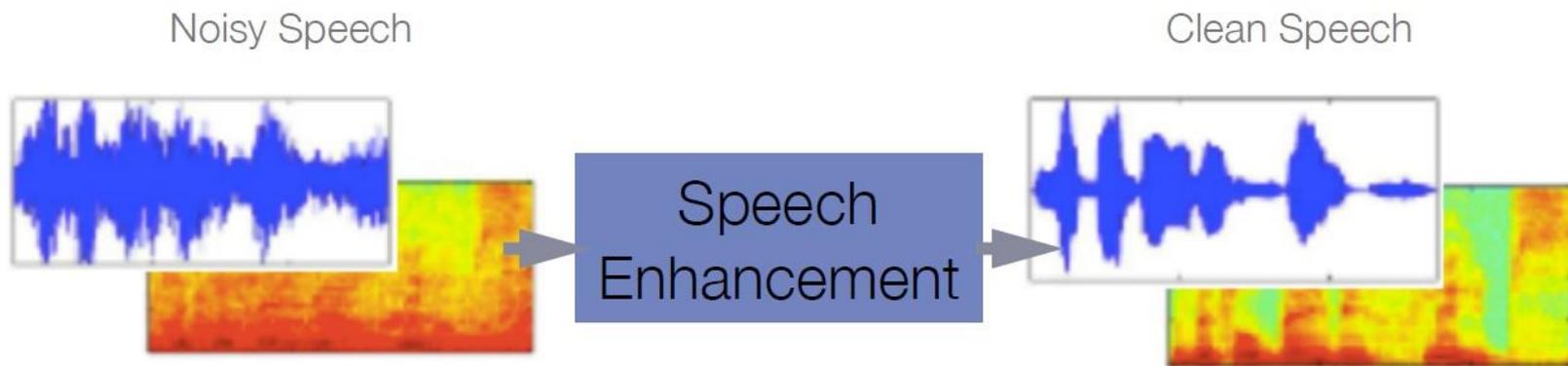**2.** Extending few-shot learning to more shots.

**Proposing**
**Novel algorithm to overcome few-shot learning limitations.**
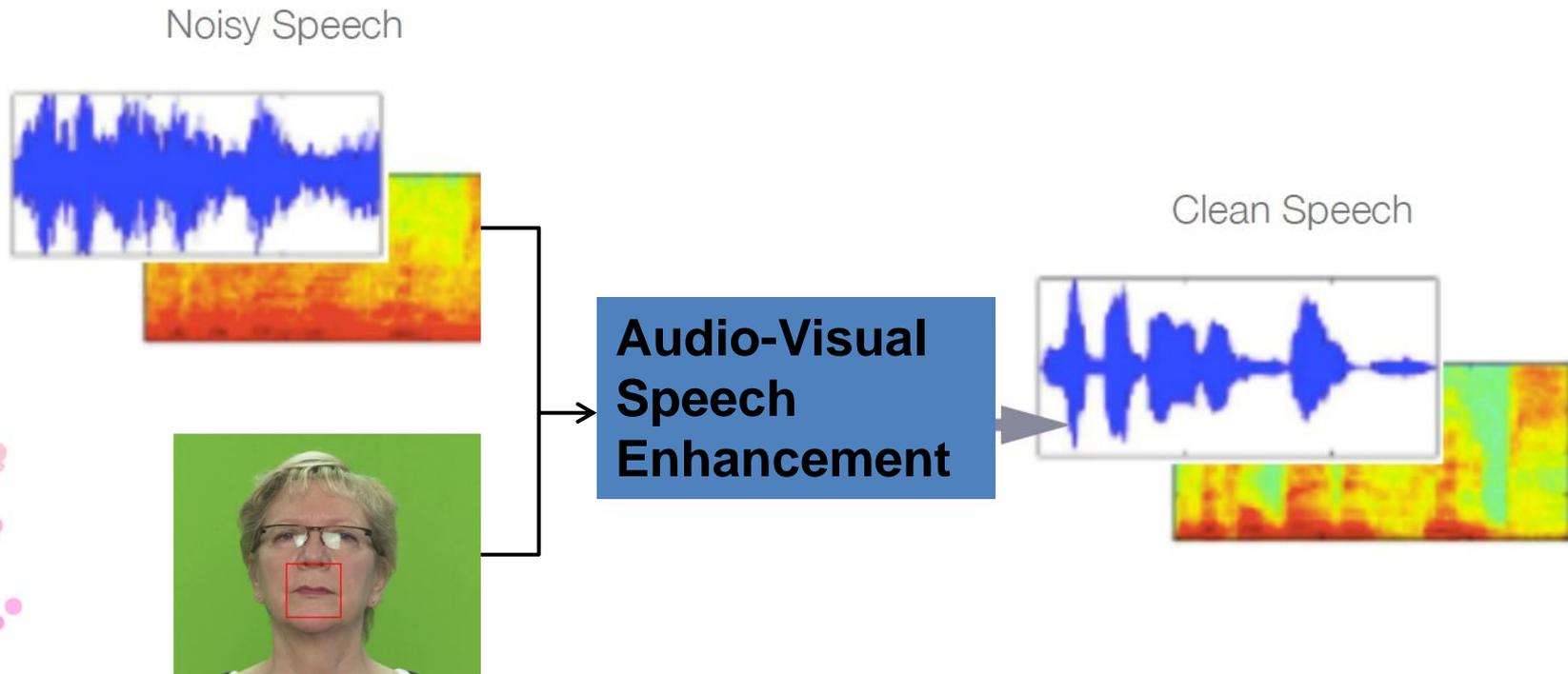**Reduce dependency on the number of shots.**

# Background & Motivation

# Speech Enhancement

Noisy Speech

Clean Speech

Speech Enhancement

# Audio-Visual Speech Enhancement

# Applications



## Real-Time

# The Problem – Speaker Dependency

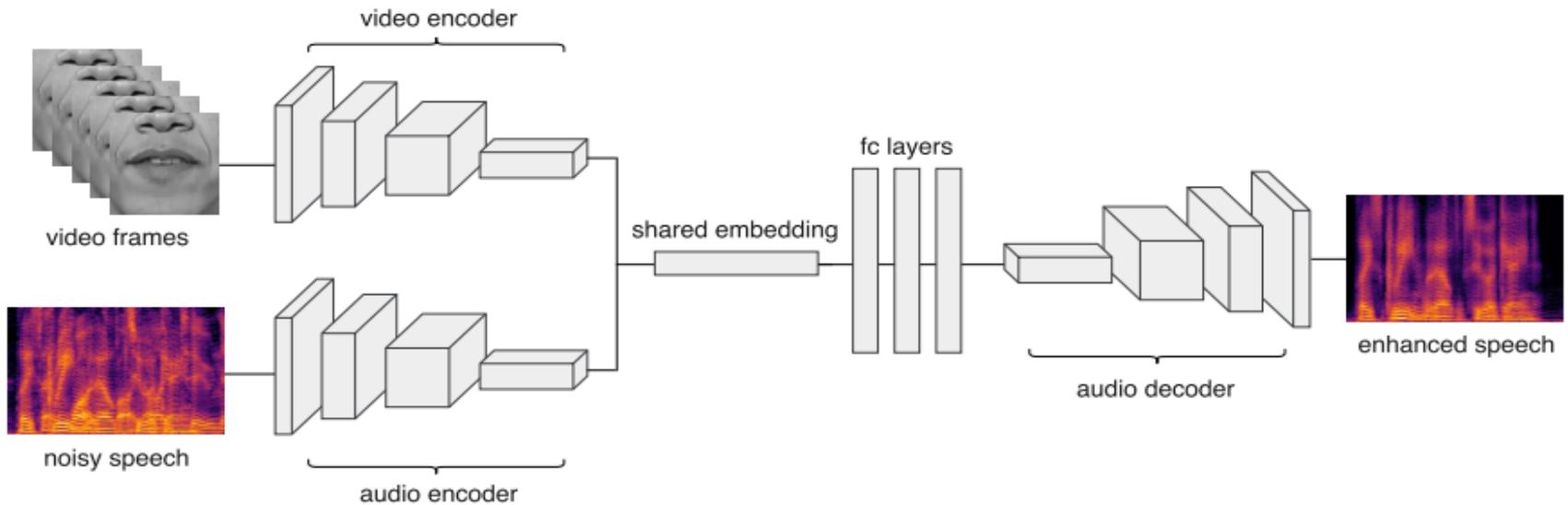Existing methods need a **previous familiarity**

classical

DL

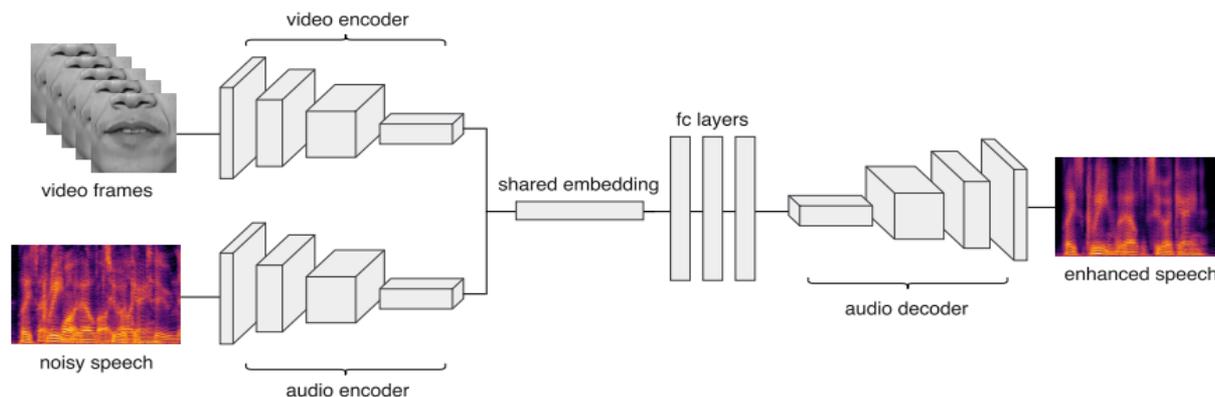Characteristics of the speaker and the noise

Large training set of the target speaker

# The Problem – Speaker Dependency



A. Gabbay, A. Shamir, and S. Peleg, "Visual speech enhancement," in Proc. Interspeech, 2018, pp. 1170–1174.
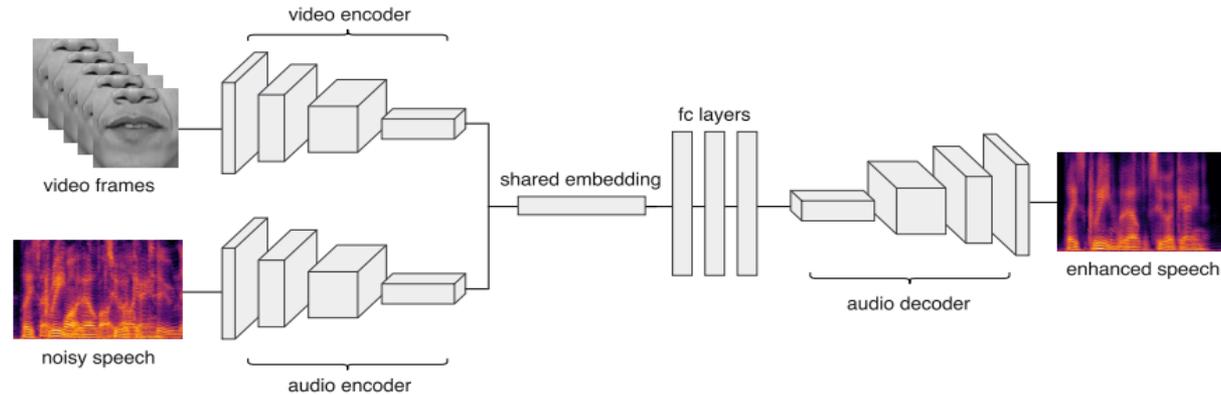
# The Problem – Speaker Dependency



| Experiment | Training | | Test | | Results (PESQ) |
|---|---|---|---|---|---|
| | **Target** | **Noise** | **Target** | **Noise** | |
| **I** | a | b | a | b | 2.86 |
| **II** | a | b | c | b | 1.989 |
| **III** | a | b | a | c | 2.105 |

A. Gabbay, A. Shamir, and S. Peleg, "Visual speech enhancement," in Proc. Interspeech, 2018, pp. 1170–1174.

# The Problem – Speaker Dependency
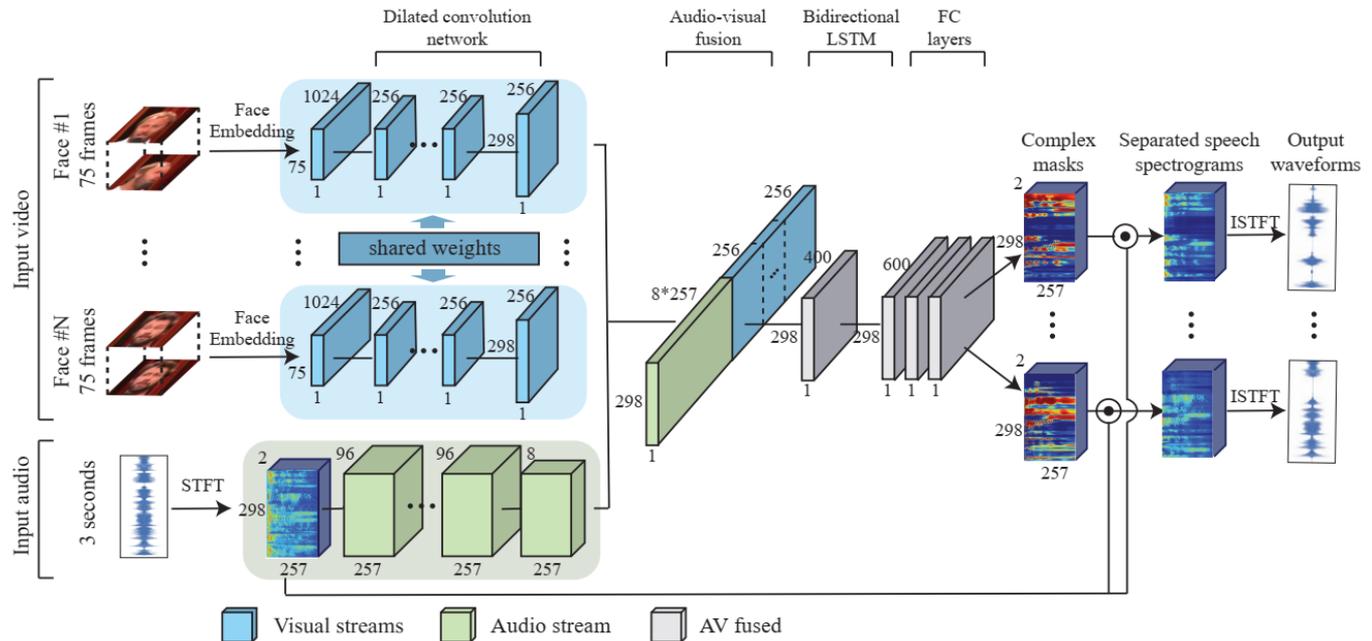


**Suffers from Speaker Dependency**

## How to Avoid it?

A. Gabbay, A. Shamir, and S. Peleg, "Visual speech enhancement," in Proc. Interspeech, 2018, pp. 1170–1174.

# Avoiding Speaker Dependency

**Very Large Dataset**



**High Computational Power**

A. Ephrat, I. Mosseri, O. Lang, T. Dekel, K. Wilson, A. Hassidim, W. T. Freeman, and M. Rubinstein, "Looking to listen at the cocktail party: A speaker-independent audio-visual model for speech separation," *arXiv preprint arXiv:1804.03619*, 2018

# Applications



our country
smart presi

0:09 / 0:34

zoom meetings

70°
Good morning!
it's 6:30 AM
dismiss

Real-Time

# **Research Contributions**

**1.** Overcoming speaker dependency for real-time mobile applications.

**Proposing**
**Fast Adaptation Speech Enhancement (FASE) model,**
**Inspired by few-shot learning methods**

**2.** Extending few-shot learning to more shots.

**Proposing**
**Novel algorithm to overcome few-shot learning limitations.**
**Reduce dependency on the number of shots.**

# Research Contributions

1. Overcoming speaker dependency for real-time mobile applications.

   **Proposing**
   **Fast Adaptation Speech Enhancement (FASE) model,**
   **Inspired by few-shot learning methods**

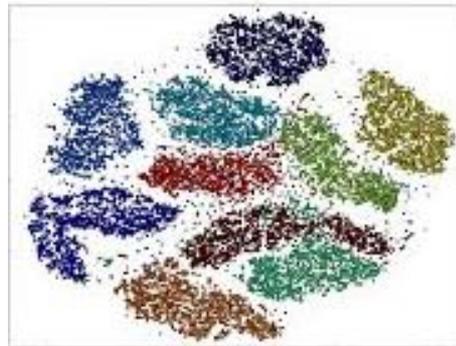2. Extending few-shot learning to more shots.

   **Proposing**
   **Novel algorithm to overcome few-shot learning limitations.**
   **Reduce dependency on the number of shots.**

# Few-Shot Learning

INPUT
(Saiga)

CONVOLUTION + RELU    POOLING    CONVOLUTION + RELU    POOLING    FLATTEN    FULLY CONNECTED    SOFTMAX

FEATURE LEARNING    CLASSIFICATION

Not enough data to learn all classes!

Goal:

- Efficiently learn *novel* categories

- Not forget the *base* categories

***n-way k-shot***     n = 1, 5
k = 1, 5

Very few compared to 600!

16

# Main Approaches

**Metric-Learning:**

Learn <u>feature representations</u> that preserve the class neighborhood structure.
→ Intra-class similarity and inter-class dissimilarity.



**Meta-Learning (Learning to Learn):**

Generate parameter updates that will optimize the classification performance of a <u>learner model on a task</u>.

| Meta-Learning [9] | task$_1$ model$_1$ | ... | task$_N$ model$_N$ | task$_{N+1}$ model$_{N+1}$ |
|---|---|---|---|---|

# Research Contributions

1. Overcoming speaker dependency for real-time mobile applications.
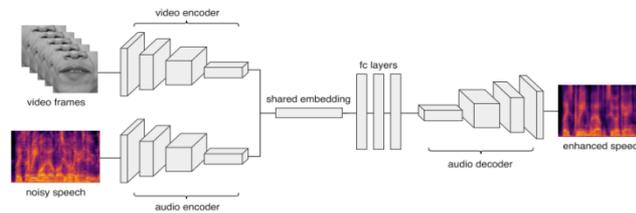
   **Proposing**
   **Fast Adaptation Speech Enhancement (FASE) model,**
   **Inspired by few-shot learning methods**

2. Extending few-shot learning to more shots.
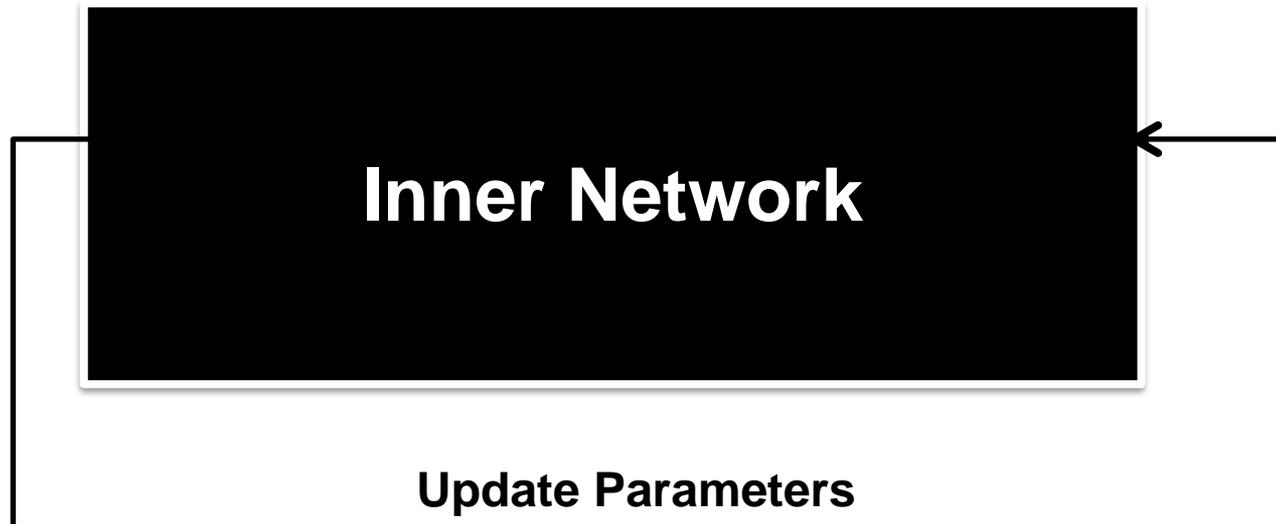
   **Proposing**
   **Novel algorithm to overcome few-shot learning limitations.**
   **Reduce dependency on the number of shots.**

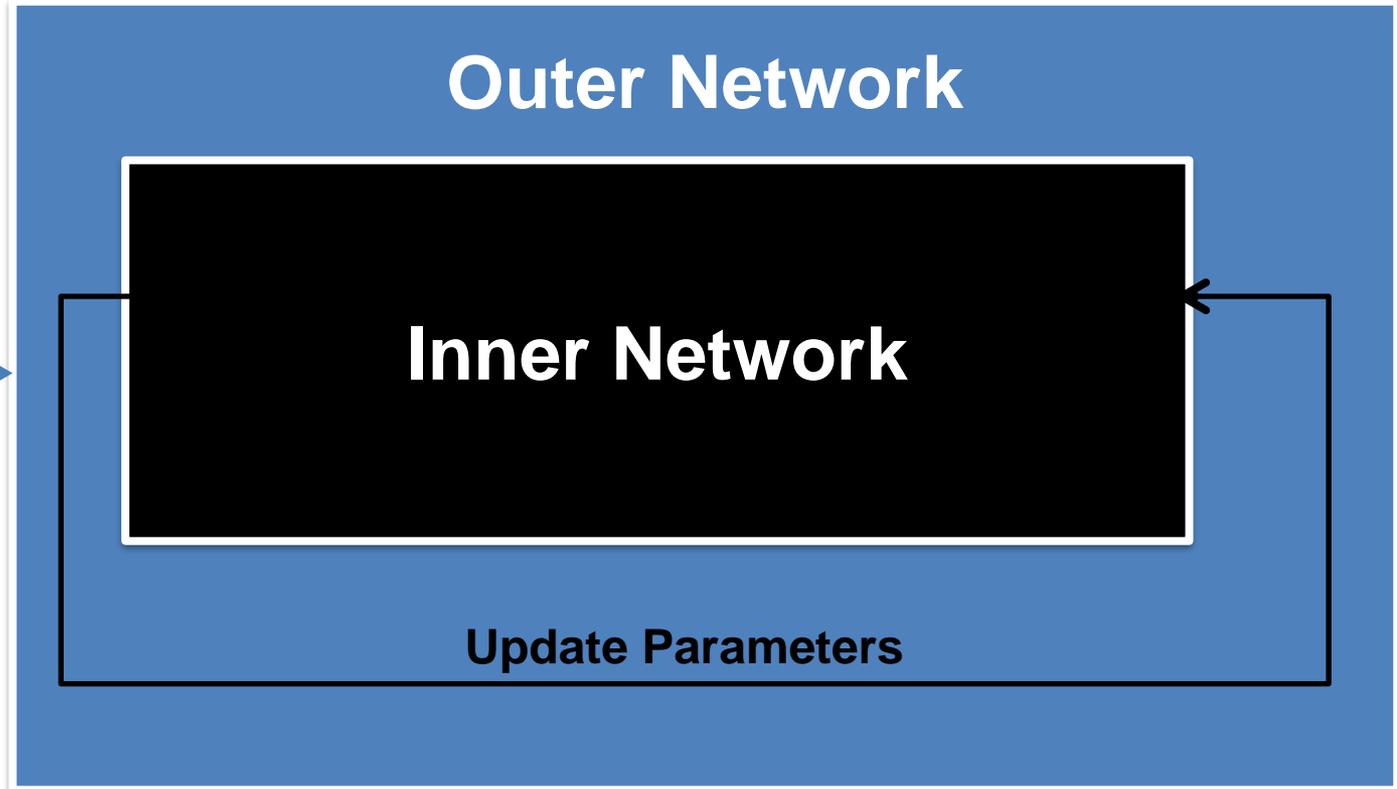# FASE: Fast Adaptation Speech Enhancement



A. Gabbay, A. Shamir, and S. Peleg, "Visual speech enhancement," in Proc. Interspeech, 2018, pp. 1170–1174.
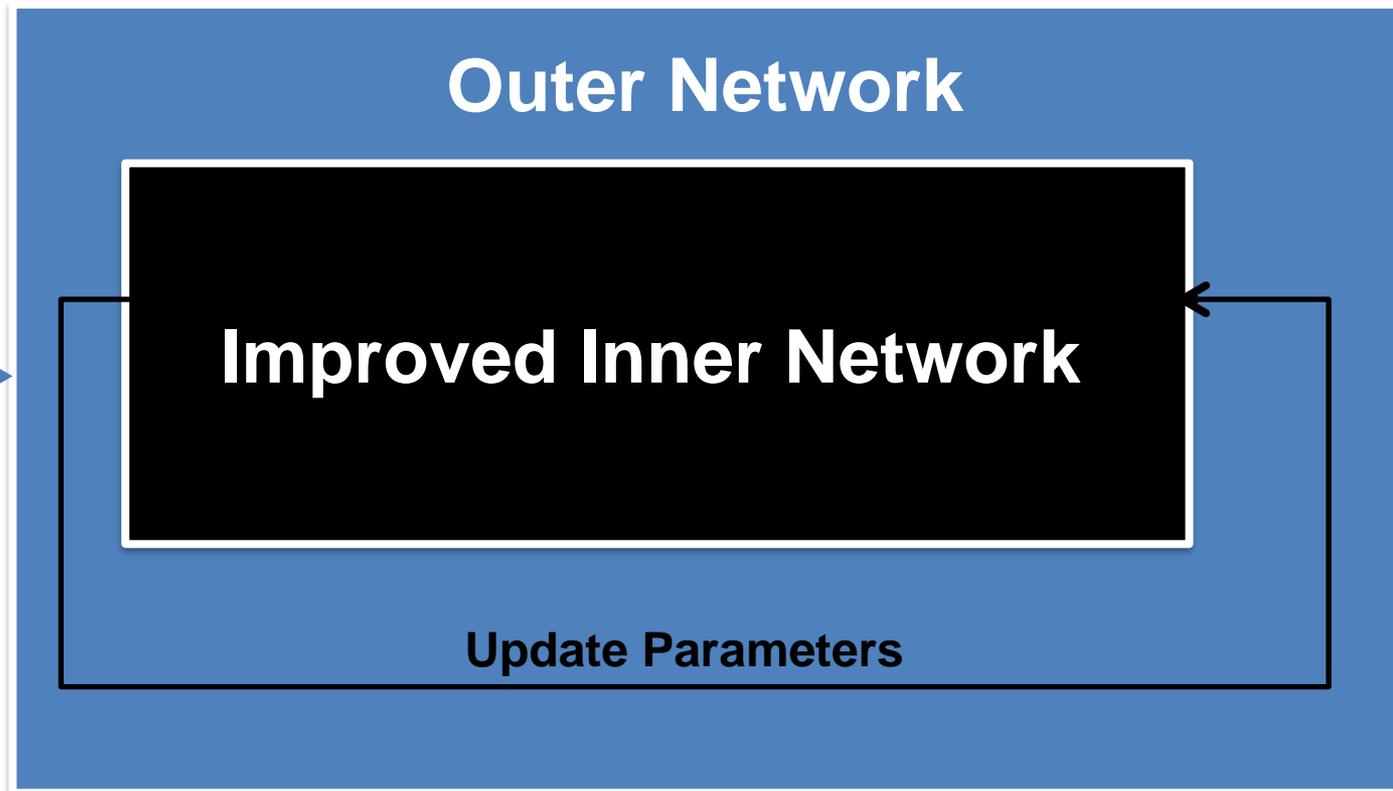
# FASE: Fast Adaptation Speech Enhancement

**Inner Network**

**Update Parameters**

# FASE: Fast Adaptation Speech Enhancement

## FASE-baseline

**Outer Network**

**Inner Network**

**Update Parameters**

Few-Shot Tasks

**Learns how to adapt to new speakers**

# FASE: Fast Adaptation Speech Enhancement

## FASE-opt



Outer Network

Improved Inner Network

Few-Shot Tasks

Update Parameters

**Faster adaptation and convergence without overfitting**

# FASE: Fast Adaptation Speech Enhancement

**FASE Training Algorithm:**

Choose number of shots *k*

Create a random pool of few-shot speech enhancement tasks.

For each few-shot task:

Train the outer network (for M epochs):

For each video:

Train the inner network (for m epochs)

Validate the inner network

Update the outer network

Validate the outer network

# Experimental Settings

**Dataset**: TCD-TIMIT

**Algorithms**: <span style="color:#8ba946">Gabbai et al.</span>
<span style="color:red">FASE-baseline</span>
<span style="color:#2e8bc0">FASE-opt</span>

**Measurements**: PESQ (**Q**uality)
STOI (**I**ntelligibility)

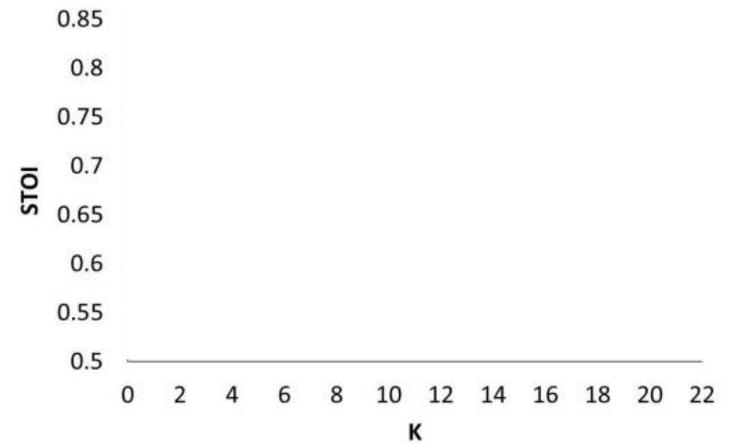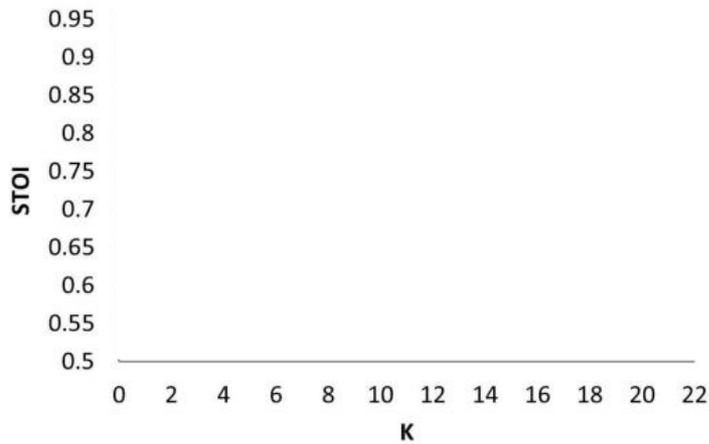**Number of Shots:** $1 \le k \le 20$

# Let's Hear It ☺

Most challenging: One-Shot Learning ($k = 1$)

**Base**

Mixture

Gabbai et al.

FASE-baseline

FASE-opt

**Novel**
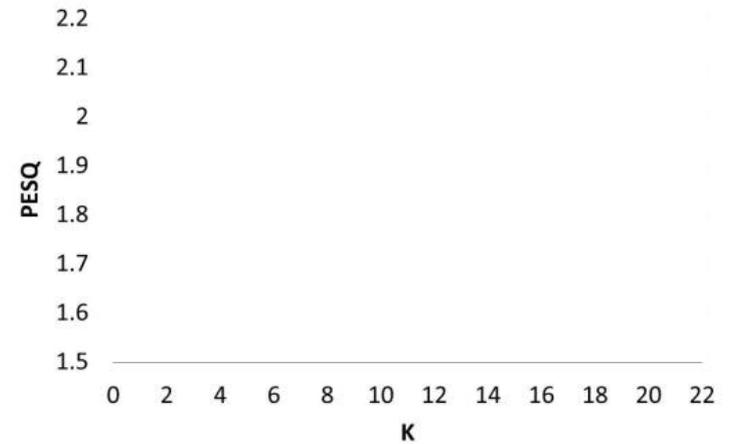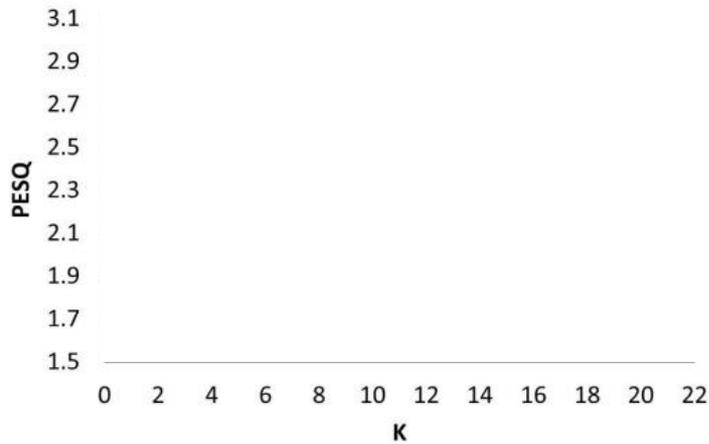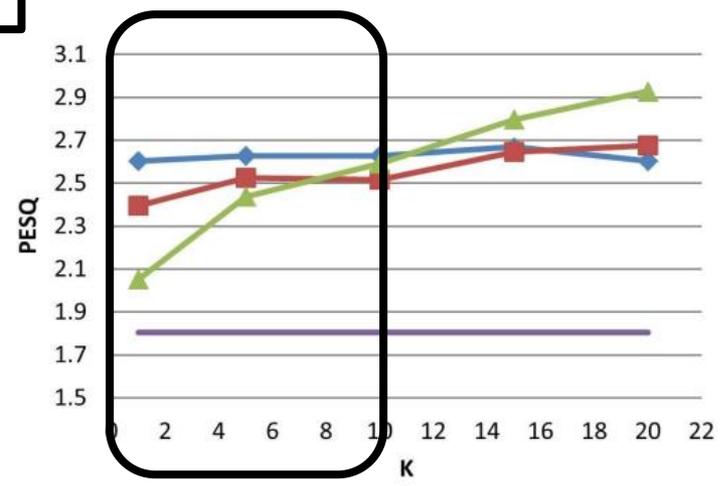
Mixture

Gabbai et al.

FASE-baseline

FASE-opt


Mixture

Gabbai et al.

FASE-baseline

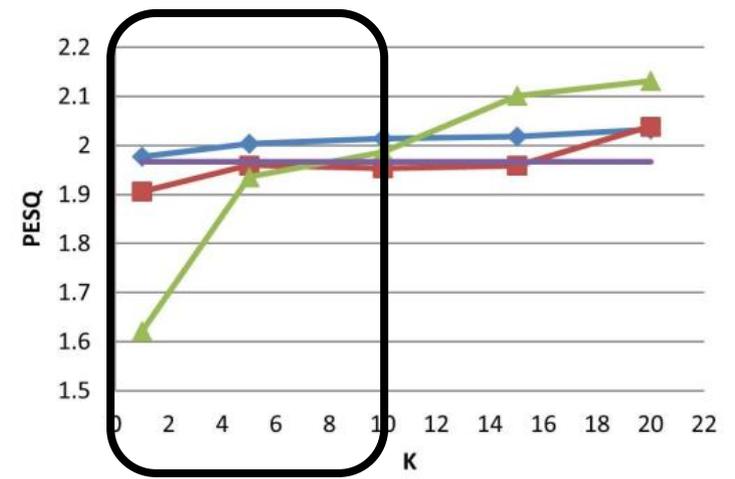FASE-opt

# FASE Model - Results



(a) Base Speakers - PESQ Score

(b) Novel Speakers - PESQ Score

(c) Base Speakers - STOI Score

(d) Novel Speakers - STOI Score

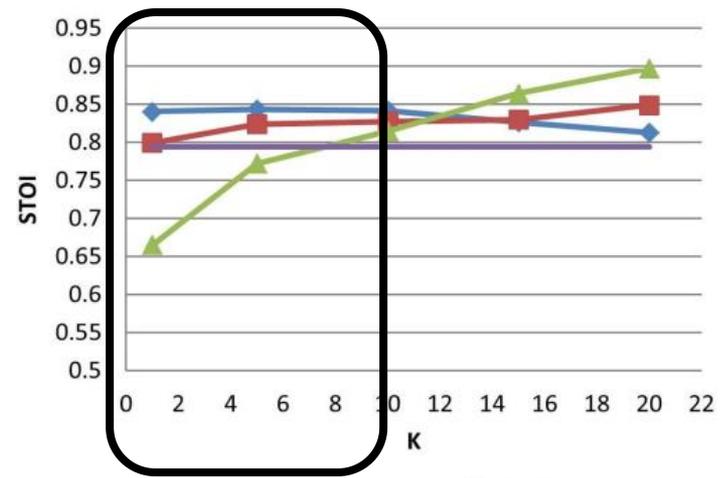FASE-opt    FASE-baseline    Gabbay et al.    Input Mixture

# FASE Model - Results



(a) Base Speakers - PESQ Score

(b) Novel Speakers - PESQ Score

(c) Base Speakers - STOI Score

(d) Novel Speakers - STOI Score

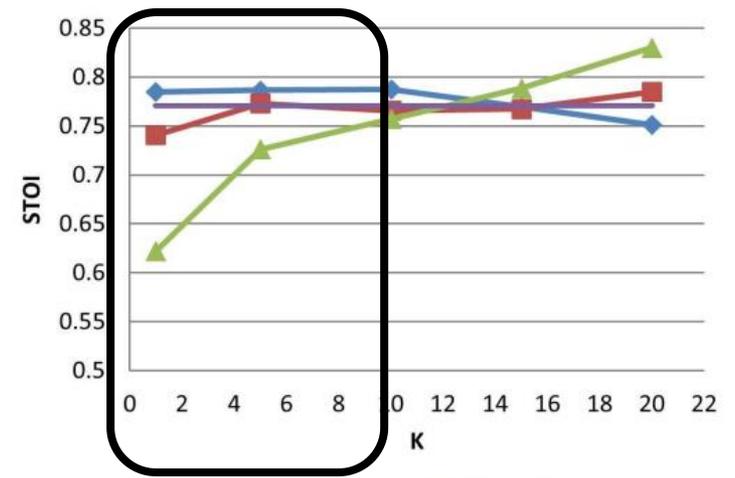FASE-opt    FASE-baseline    Gabbay et al.    Input Mixture

27

# FASE Model - Results



(a) Base Speakers - PESQ Score

(b) Novel Speakers - PESQ Score

(c) Base Speakers - STOI Score

(d) Novel Speakers - STOI Score
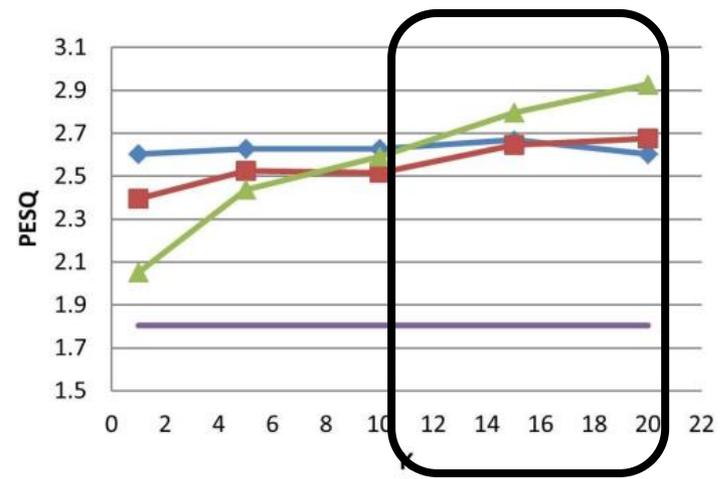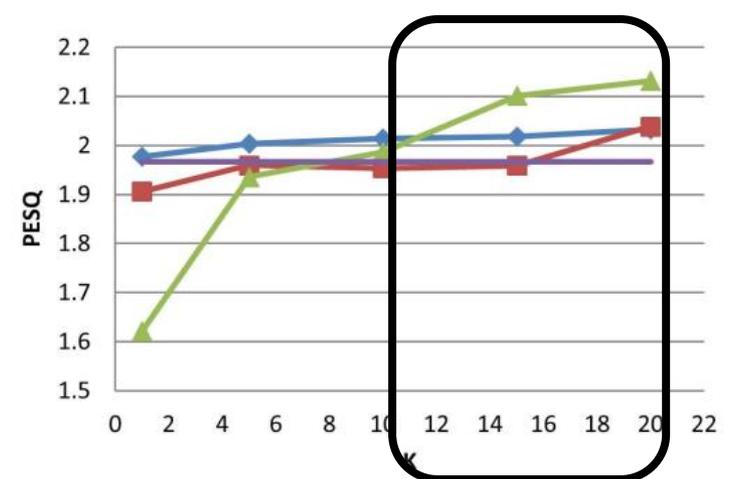
FASE-opt    FASE-baseline    Gabbay et al.    Input Mixture

28

# Let's Hear It

Largest gap: $k = 20$

**<u>Base</u>**

<u>Mixture</u>

<u>Gabbai et al.</u>

<u>FASE-baseline</u>

<u>FASE-opt</u>

**<u>Novel</u>**

<u>Mixture</u>

<u>Gabbai et al.</u>

<u>FASE-baseline</u>

<u>FASE-opt</u>

# FASE Model - Summary

**Fast Adaptation Speech Enhancement**
**Avoids Speaker dependency**

"Standard"
Few-Shot
Learning

Cases of
More Shots

✓ Quality (PESQ)
✓ Intelligibility (STOI)
✓ Less Computational Power

Let's investigate...

# Research Contributions

**1.** Overcoming speaker dependency for real-time mobile applications.

**Proposing**
**Fast Adaptation Speech Enhancement (FASE) model,**
**Inspired by few-shot learning methods**

**2.** Extending few-shot learning to more shots.

**Proposing**
**Novel algorithm to overcome few-shot learning limitations.**
**Reduce dependency on the number of shots.**

# **Research Contributions**

1. Overcoming speaker dependency for real-time mobile applications.

   **Proposing**
   **Fast Adaptation Speech Enhancement (FASE) model,**
   **Inspired by few-shot learning methods**

2. Extending few-shot learning to more shots.

   **Proposing**
   **Novel algorithm to overcome few-shot learning limitations.**
   **Reduce dependency on the number of shots.**

# Research Contributions

1. Overcoming speaker dependency for real-time mobile applications.

   **Proposing**
   **Fast Adaptation Speech Enhancement (FASE) model,**
   **Inspired by few-shot learning methods**

2. Extending few-shot learning to more shots.

   **Proposing**
   **Novel algorithm to overcome few-shot learning limitations.**
   **Reduce dependency on the number of shots.**

# Few-Shot Learning Limitations

- Customized for $k = 1$, $k = 5$
- In reality – No guarantee
- From few-shots to more shots:
  - Exploring limitations
  - Proposed algorithm
  - Results (proof-of-concept)
  - Conclusions & Future work
- Gidaris et al. *

* S. Gidaris and N. Komodakis, "Dynamic few-shot visual learning without forgetting," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2018, pp. 4367–4375.
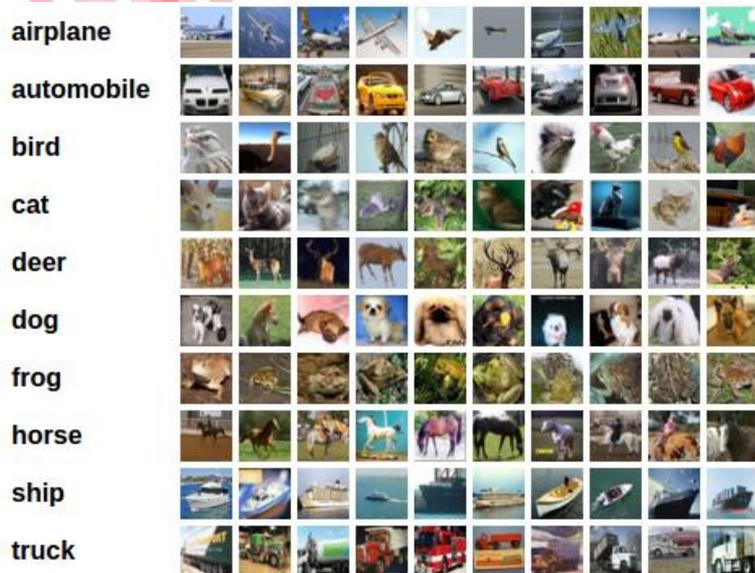
# Experimental Settings

**Dataset**: mini-ImageNet

**Algorithms**:  Gidaris et al.
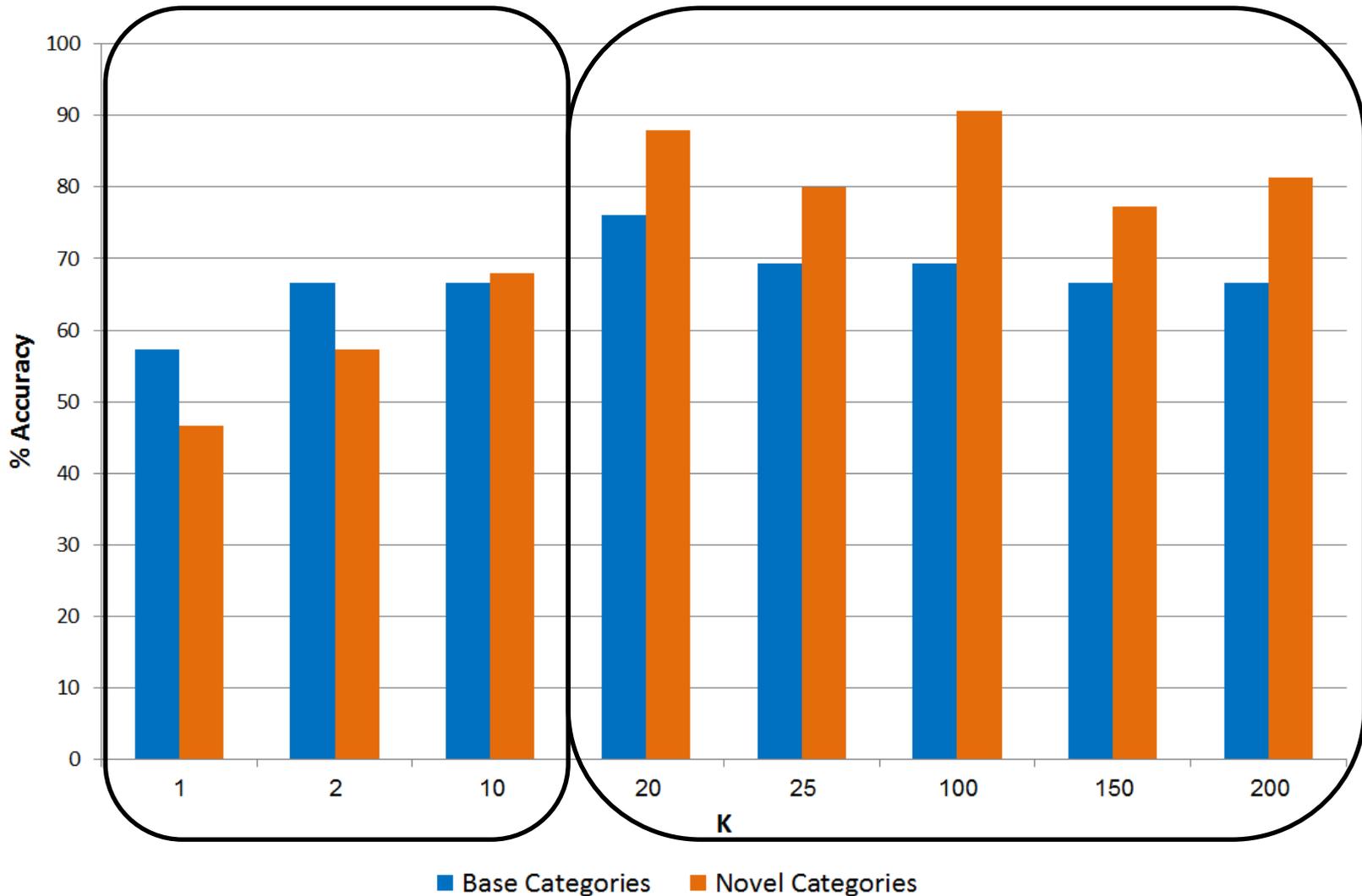                 Later – Our solution

**Measurements**:  % Accuracy

**Number of Shots:**  $1 \leq k \leq 200$

# Few-Shot Learning Limitations

STD along all cases: 9.25



36

# From Few-Shot to More Shots

## Proposed algorithm:

Train the feature extractor of Gidaris et al.
Divide the base feature vectors evenly into $N$ spaces.
Continue as Gidaris et al.
For each image:
    Calculate matching in all $N$ spaces
    Choose the best category by weighted calculation
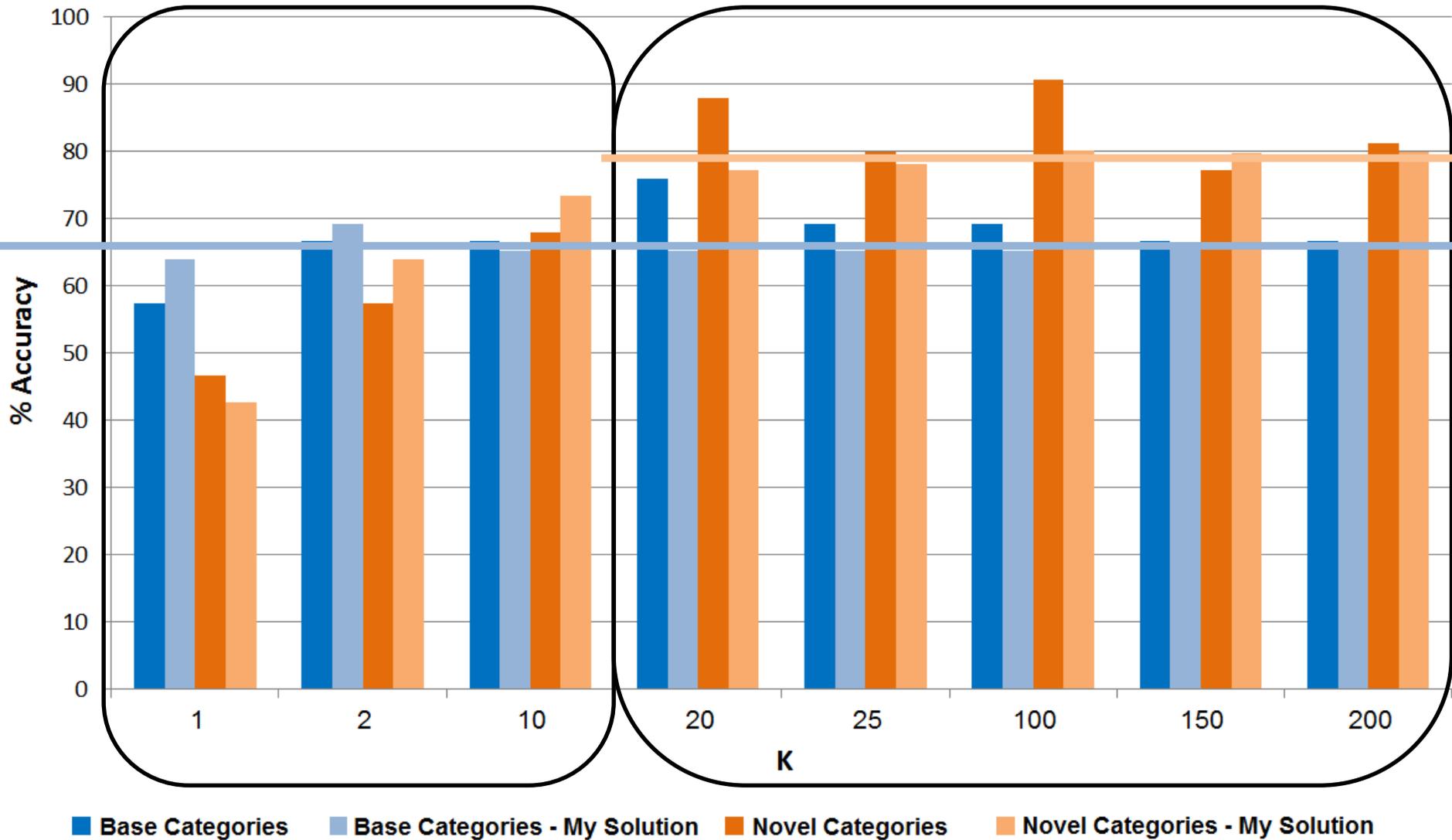
## Proof-of-concept ☺

Dividing method = Random
Training iterations = 20

# Results

**Base Categories**    **Base Categories - My Solution**    **Novel Categories**    **Novel Categories - My Solution**

# More Shots - Summary

Results:

✓ Smaller base-novel trade-off

✓ More shots: Improved stability

✓ Less dependency on $k$

✓ **The method is general**
   Can be used in other few-shot tasks

Future work:

- Optimize training.
- Explore the separation method into spaces.
- Find the ideal number of spaces.

# **Research Contributions**

1. Overcoming speaker dependency for real-time mobile applications.

   ✓ Quality (PESQ)
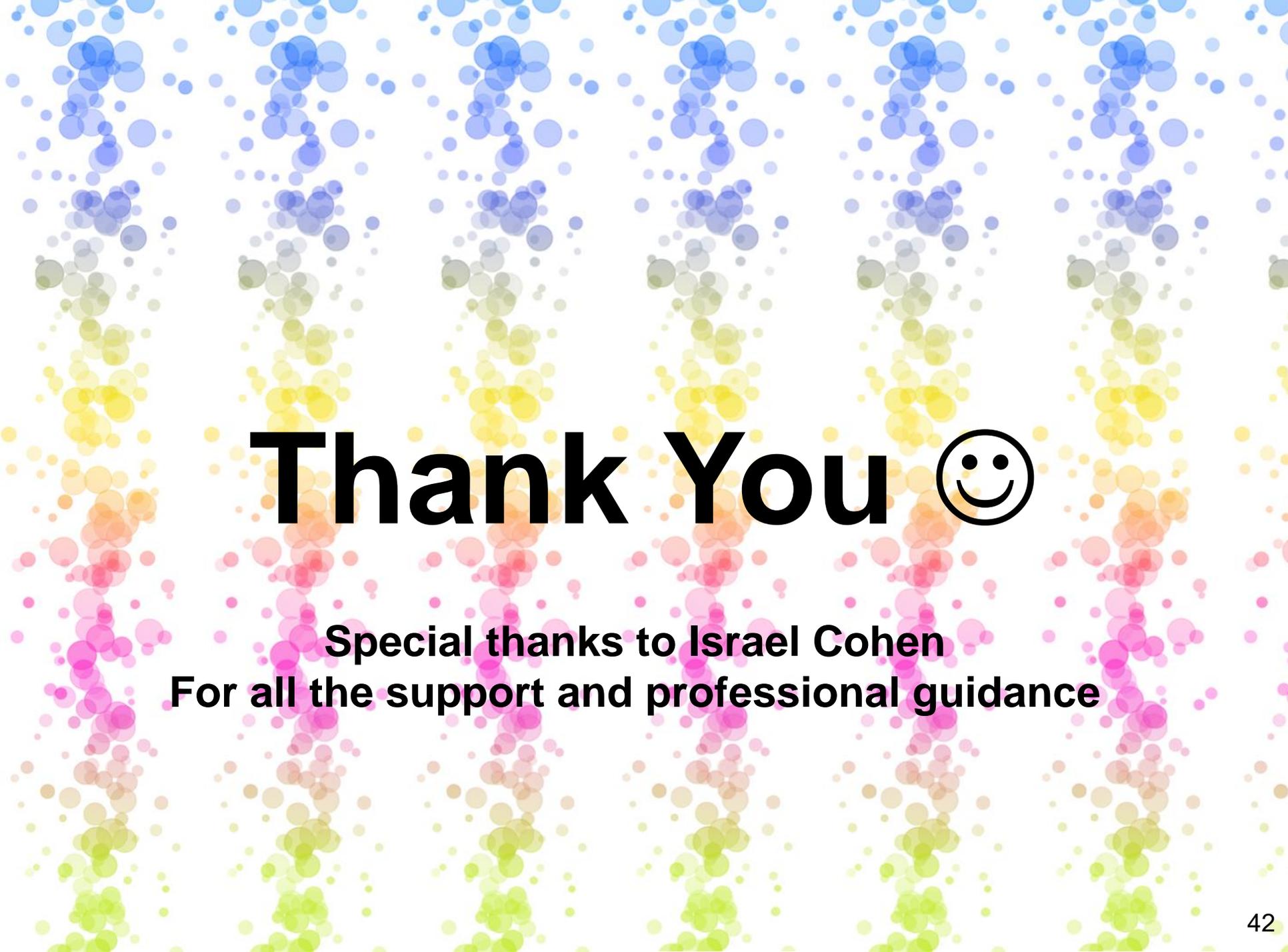   ✓ Intelligibility (STOI)
   ✓ Less computational power

2. Extending few-shot learning to more shots.

   ✓ Improved stability
   ✓ Smaller base-novel trade-off
   ✓ General method

# Future Research

- Improving few-shot learning model for more shots.
  → Training, space separation, number of spaces

- Improving results of FASE model with few shots.
  → Larger dataset, training configurations

- Improving results of FASE model with more shots.
  → Deploying the proposed **general solution**
     in audio-visual speech enhancement
  → Learn a **filter** rather than the spectrogram itself.
     (inability to create sounds that do not exist in
     the original soundtrack)

# Thank You ☺

**Special thanks to Israel Cohen**
**For all the support and professional guidance**