

# A multi-objective perspective for the Once-for-All neural architecture search framework

MSc Defense

Rafael Claro Ito

28/02/2023



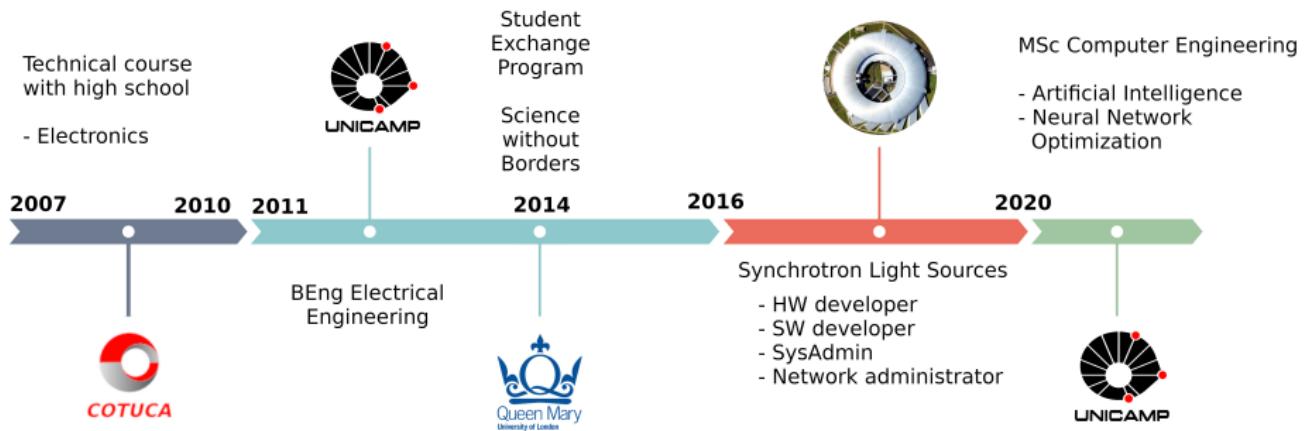
# Outline

- 1 Introduction
- 2 AutoML
- 3 Training stage
- 4 Search stage
  - Random search
  - OFA search
  - OFA<sup>2</sup> search
  - Comparison
- 5 Ensemble
  - Manual sampling
  - Multi-objective Optimization
- 6 Conclusion
- 7 References

# Agenda

- 1 Introduction
- 2 AutoML
- 3 Training stage
- 4 Search stage
  - Random search
  - OFA search
  - OFA<sup>2</sup> search
  - Comparison
- 5 Ensemble
  - Manual sampling
  - Multi-objective Optimization
- 6 Conclusion
- 7 References

# whoami?



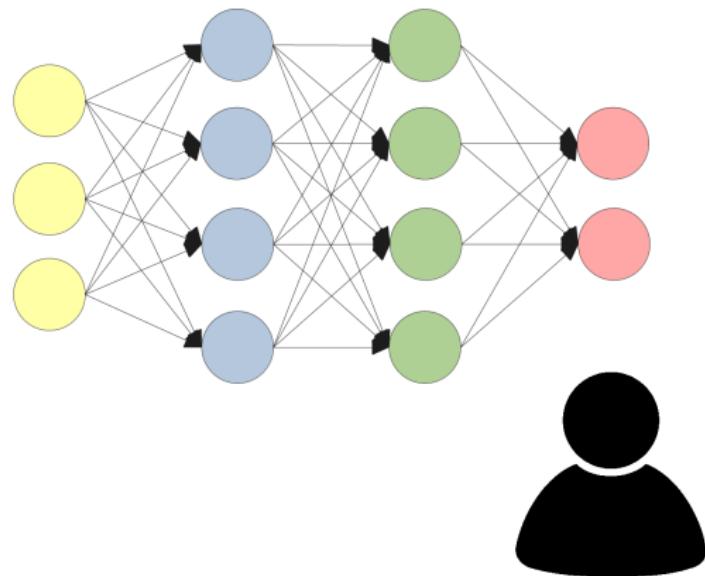
# Acknowledgements

Part of the results presented in this work were obtained through the project "Hub of Artificial Intelligence in Health and Wellbeing – Viva Bem", funded by Samsung Eletrônica da Amazônia Ltda.

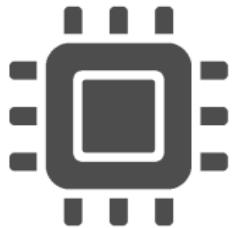


Hub of Artificial Intelligence for Health and Well-being.

# One deployment scenario

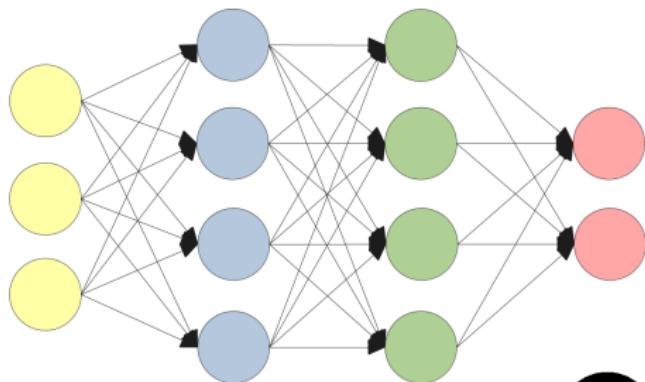


# One deployment scenario

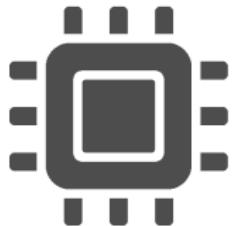


Hardware restriction

Metric	Example
memory	12 GB
latency	50 ms
FLOPS	1 G
MACs	600 M

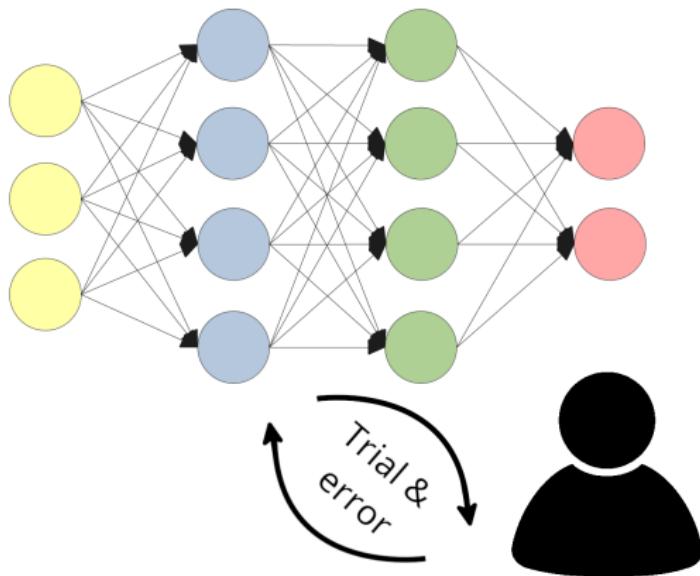


# One deployment scenario

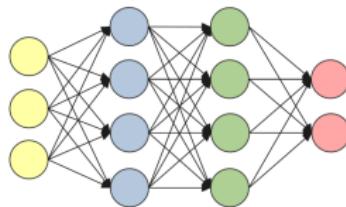
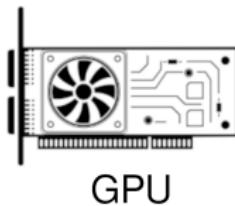


Hardware restriction

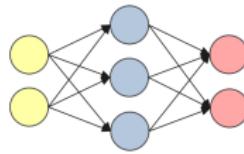
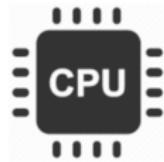
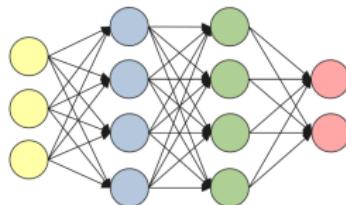
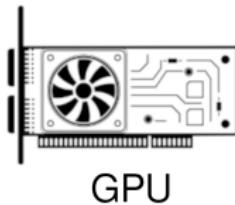
Metric	Example
memory	12 GB
latency	50 ms
FLOPS	1 G
MACs	600 M



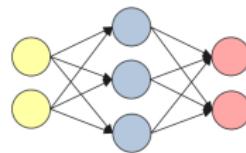
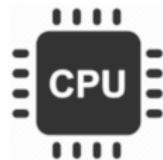
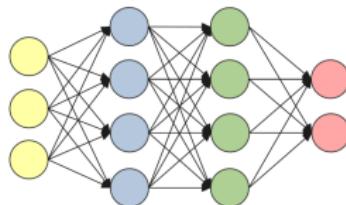
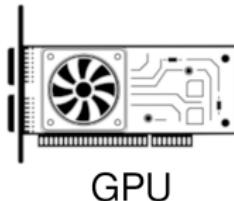
# Many deployment scenario



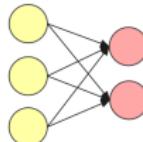
# Many deployment scenario



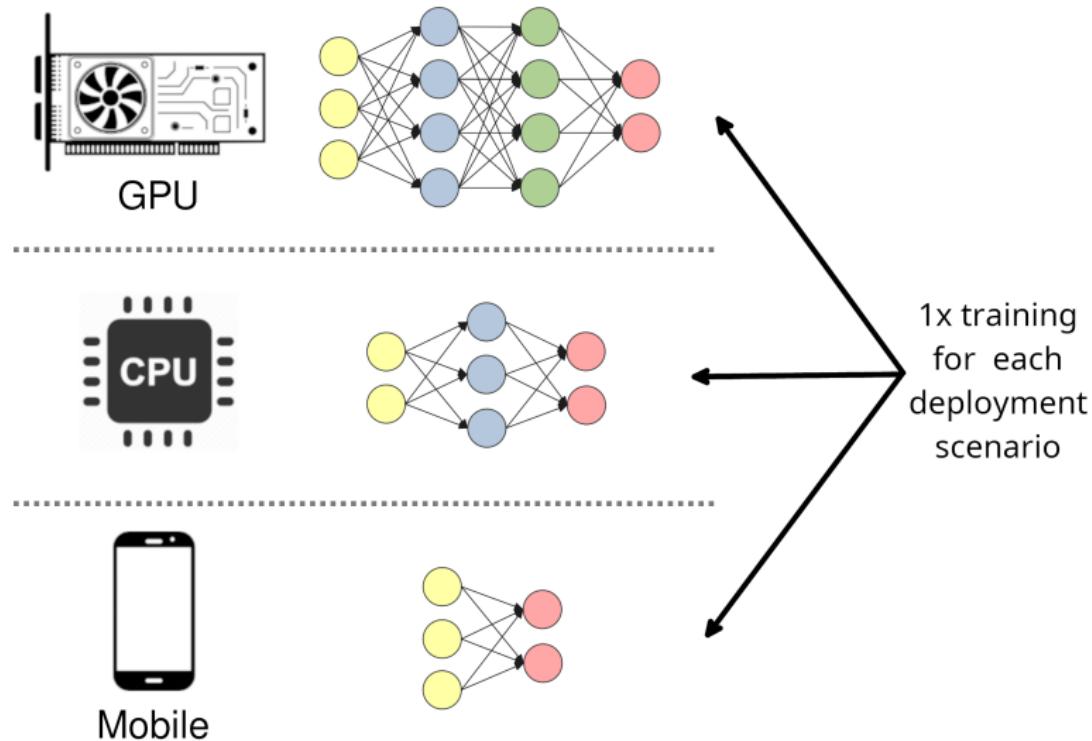
# Many deployment scenario



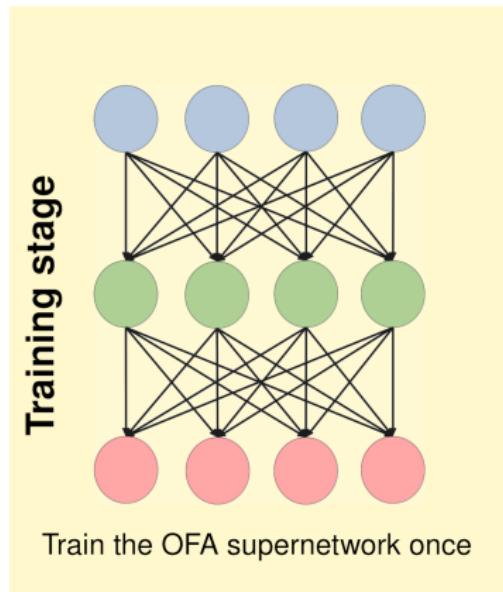
Mobile



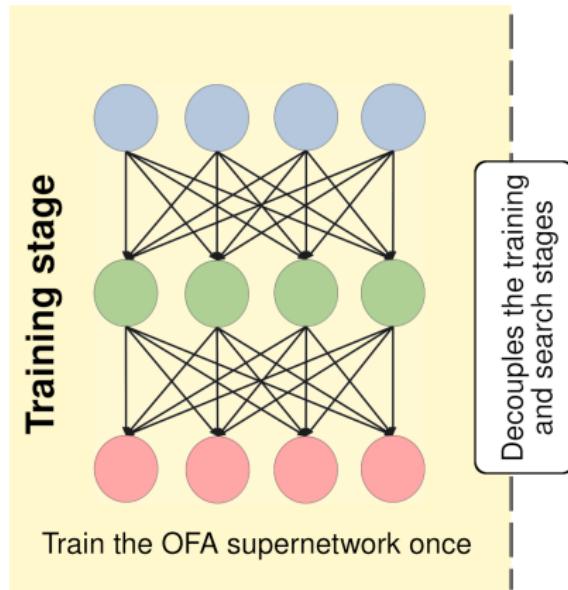
# Many deployment scenario



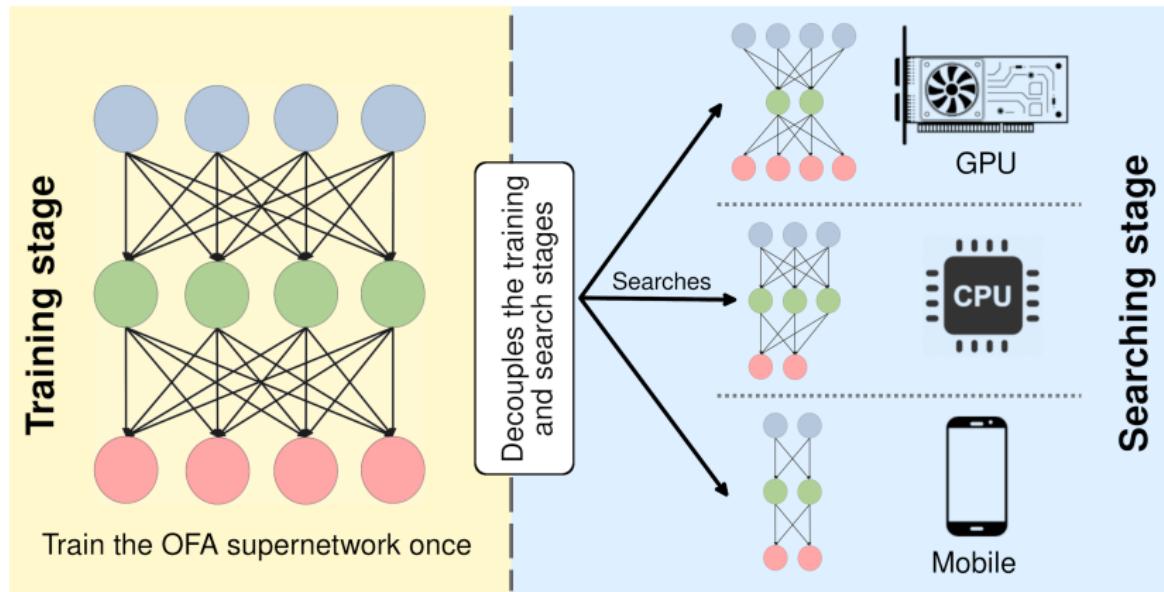
# Once-for-All framework



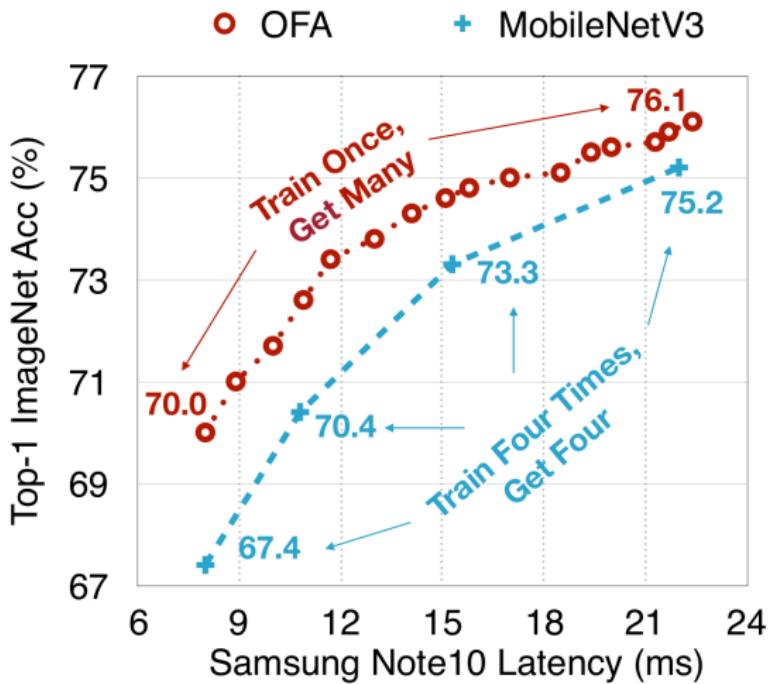
# Once-for-All framework



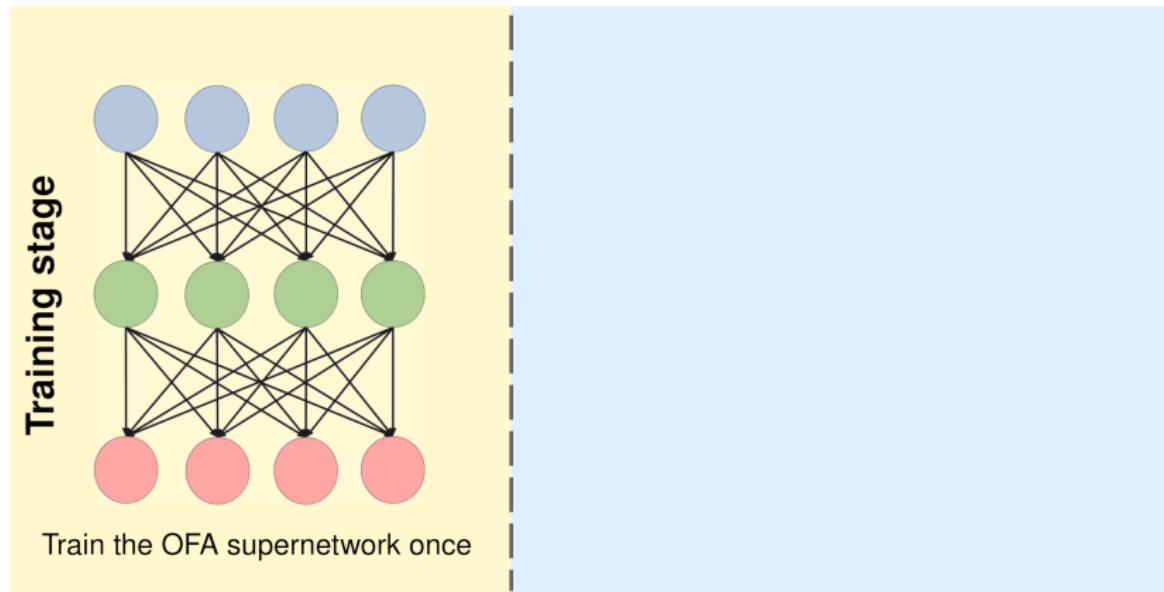
# Once-for-All framework



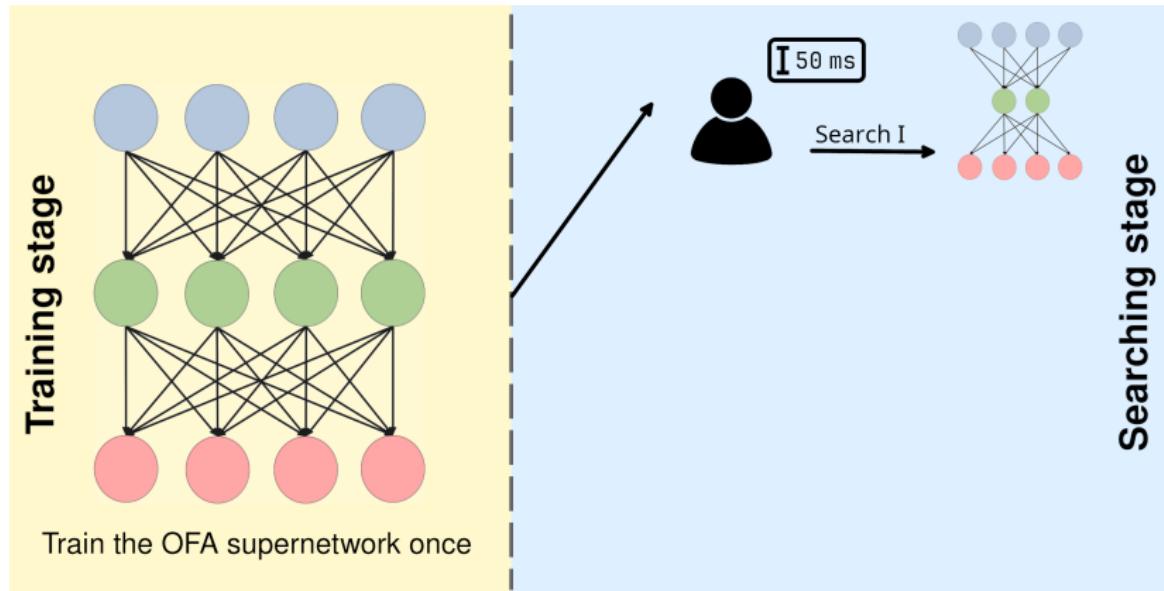
# Once-for-All framework



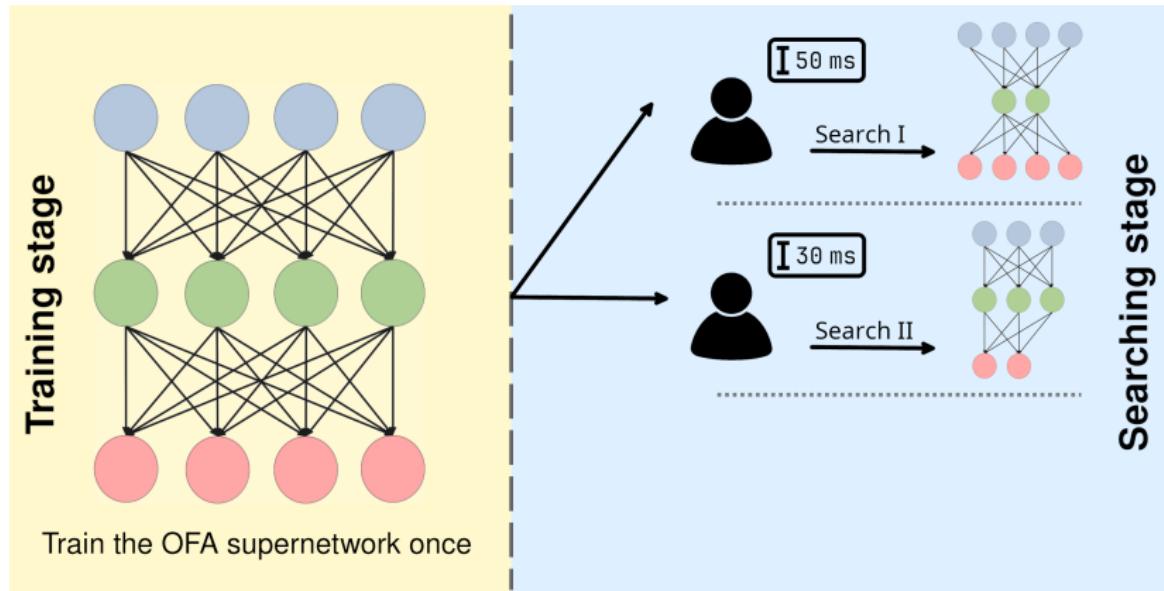
# Once-for-All framework



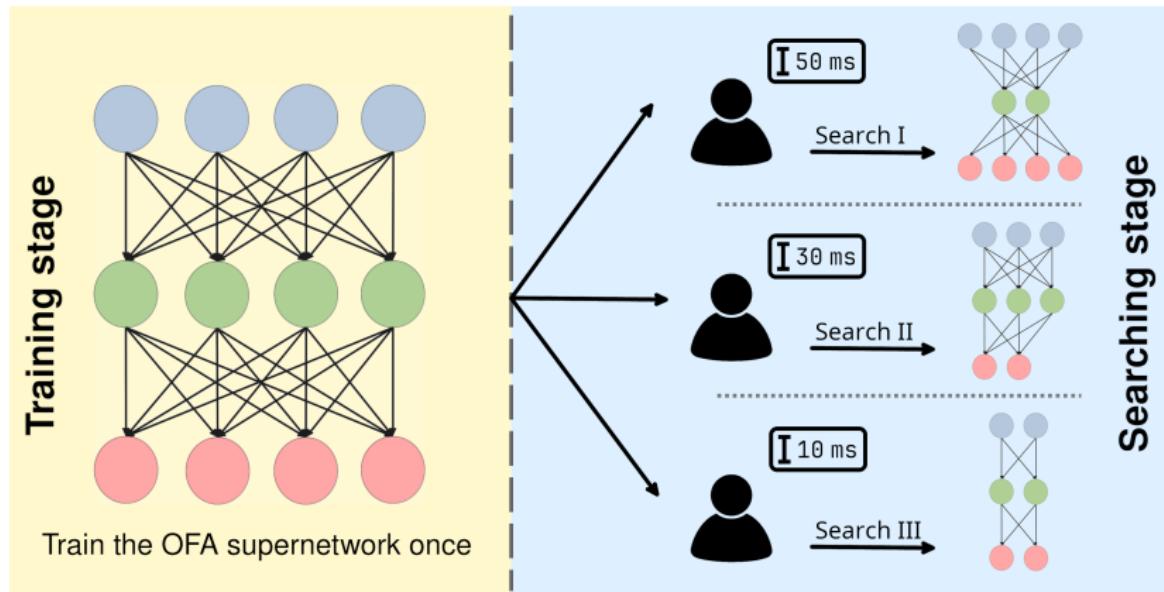
# Once-for-All framework



# Once-for-All framework



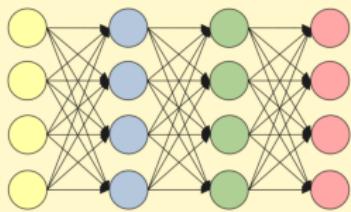
# Once-for-All framework



Train the OFA supernet once

# OFA<sup>2</sup> framework

Once-for-All

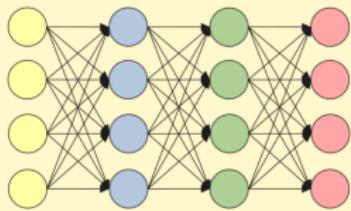


Train the OFA supernet once

**Training stage**

# OFA<sup>2</sup> framework

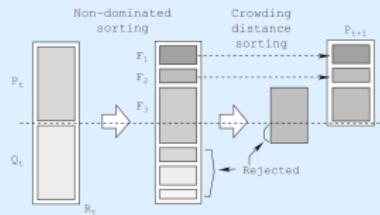
## Once-for-All



Train the OFA supernet once

## Training stage

## EMOA

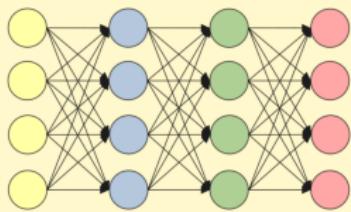


Search once with a MOO formulation

## Search stage

# OFA<sup>2</sup> framework

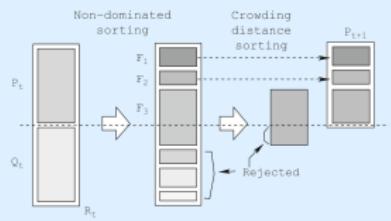
## Once-for-All



Train the OFA supernet once

## Training stage

## EMOA

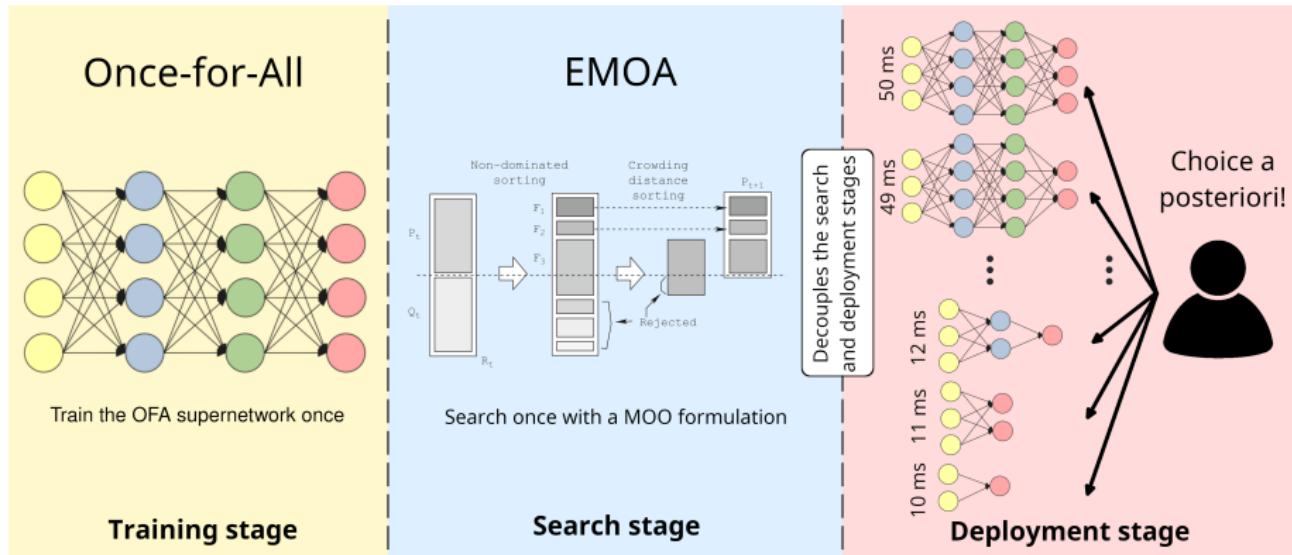


Decouples the search  
and deployment stages

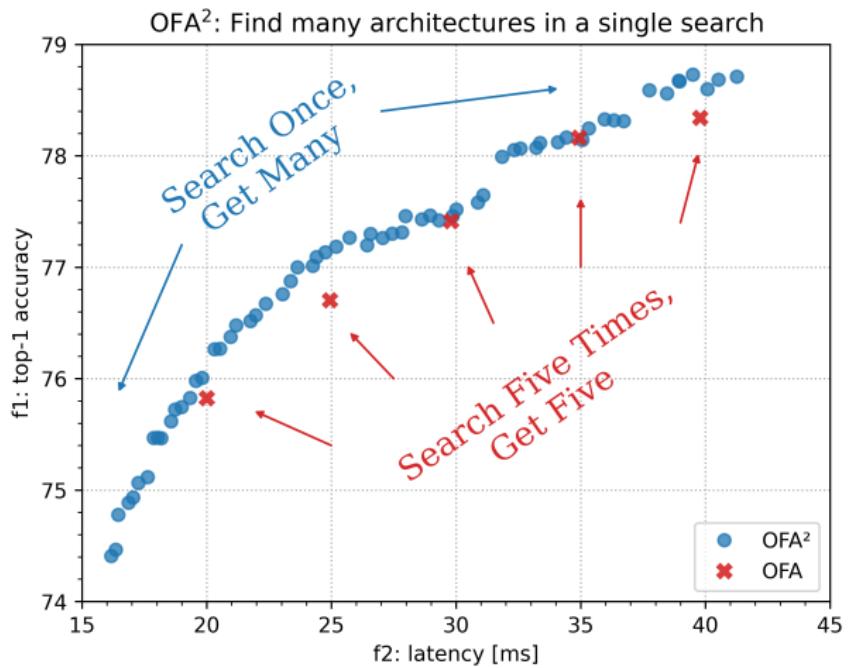
Search once with a MOO formulation

## Search stage

# OFA<sup>2</sup> framework



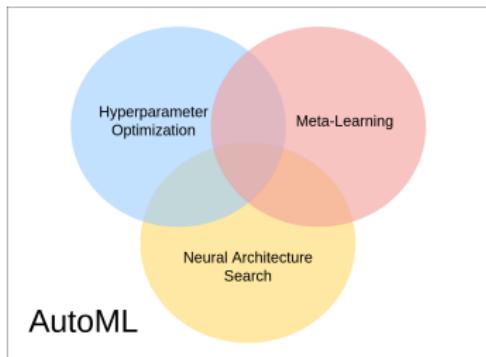
# OFA<sup>2</sup> framework



# Agenda

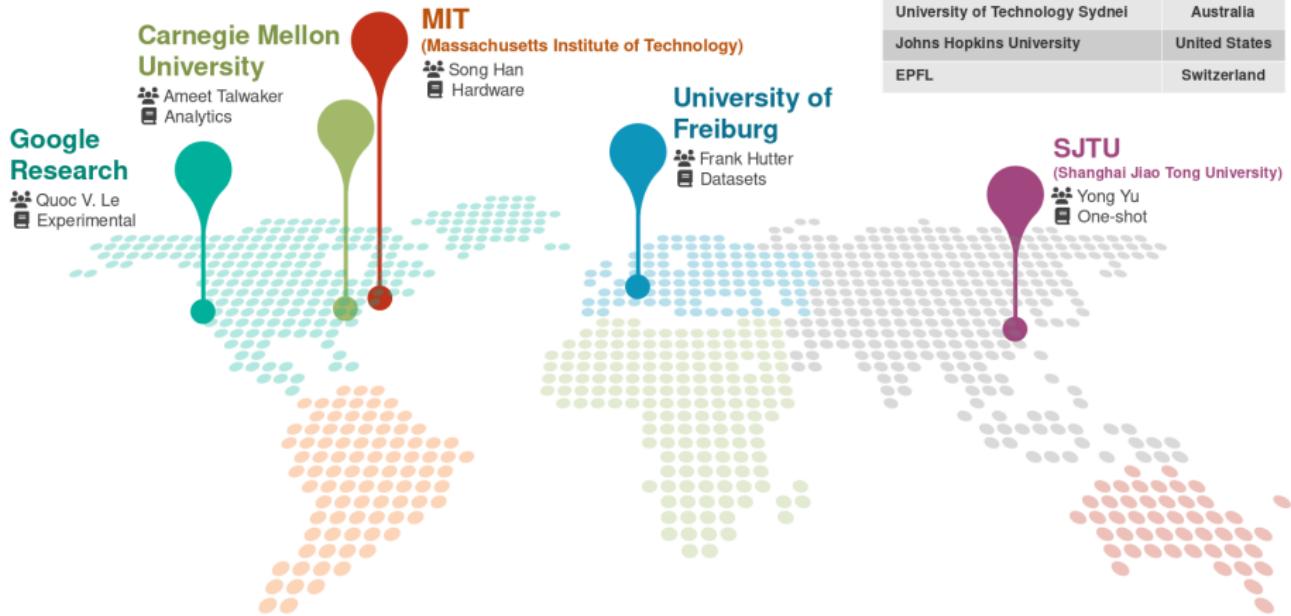
- 1 Introduction
- 2 AutoML
- 3 Training stage
- 4 Search stage
  - Random search
  - OFA search
  - OFA<sup>2</sup> search
  - Comparison
- 5 Ensemble
  - Manual sampling
  - Multi-objective Optimization
- 6 Conclusion
- 7 References

# Research areas

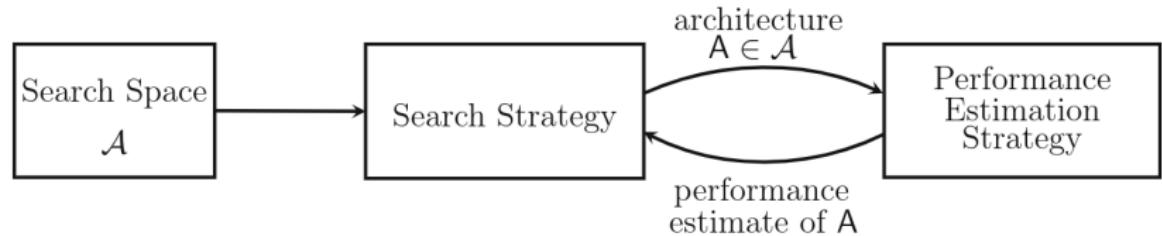


- **Hyperparameter Optimization (HPO)**
  - Tuning hyperparameters
  - Bayesian optimization
  - e.g.: Grid search
- **Meta-learning**
  - Reuse previous knowledge
  - e.g.: Transfer learning
  - e.g.: Few-shot learning
- **Neural Architecture Search (NAS)**
  - Search the network automatically

# Research groups



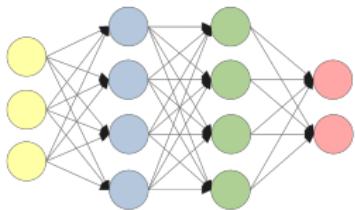
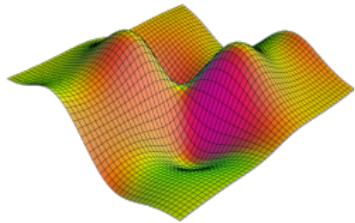
# NAS framework



(Elsken, Metzen, and Hutter 2019)

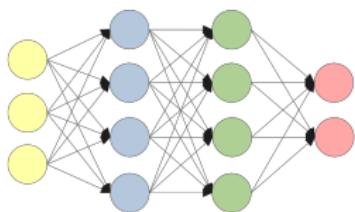
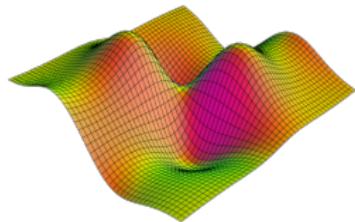
# NAS framework

Search Space

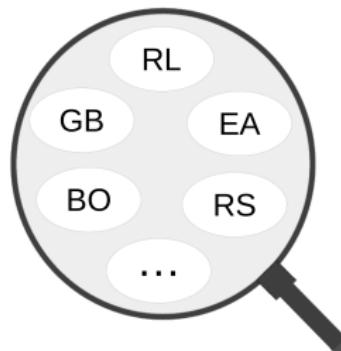


# NAS framework

Search Space

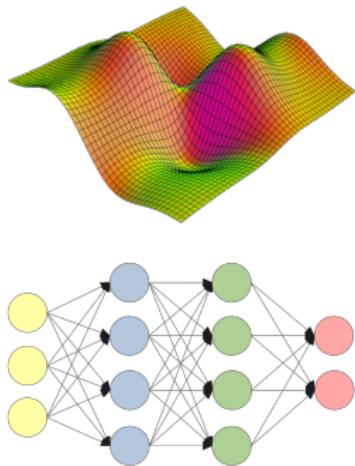


Search Strategy

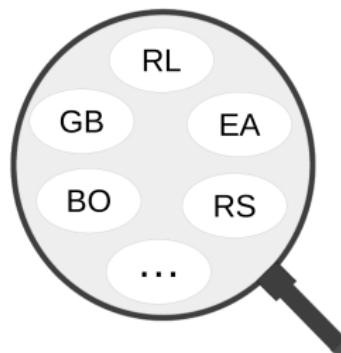


# NAS framework

## Search Space



## Search Strategy



## Search Evaluation

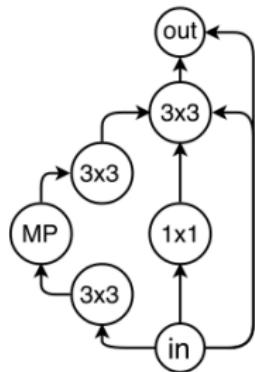
- Full training
- Partial training
- Weight-sharing
- Network morphism
- Hypernetworks
- Performance prediction

# Search Space

## Discrete Search Space:

- Vertices/Nodes: operations
- Edges: connections

Example:



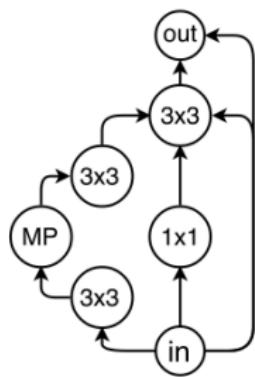
(Ying et al. 2019)

# Search Space

## Discrete Search Space:

- Vertices/Nodes: operations
- Edges: connections

Example:

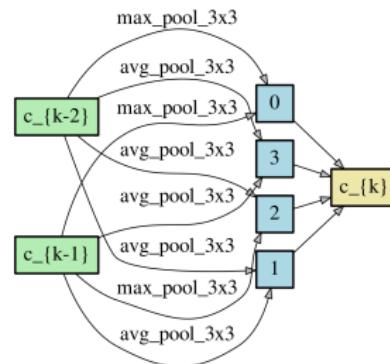


(Ying et al. 2019)

## Continuous Search Space:

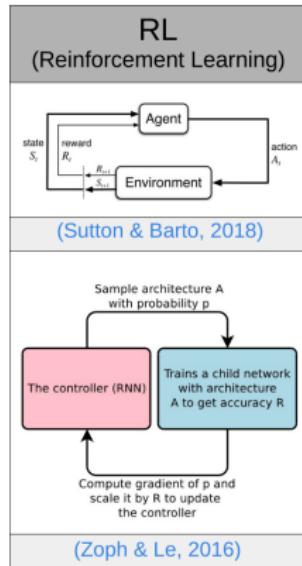
- Vertices/Nodes: latent representation (eg: feature map in conv nets)
- Edges: operations

Example:

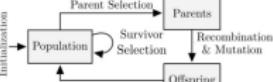
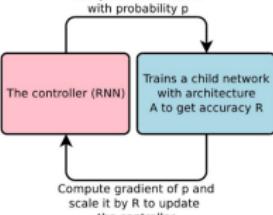
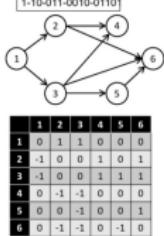


(Green et al. 2019)

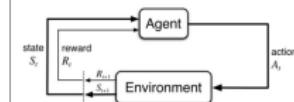
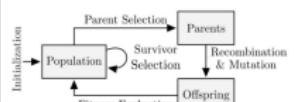
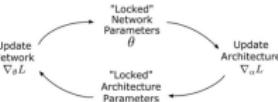
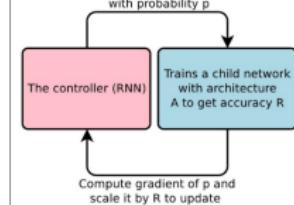
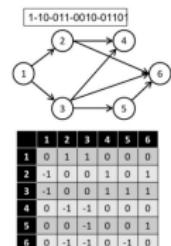
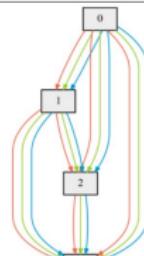
# Search Strategy



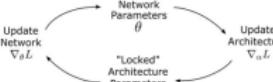
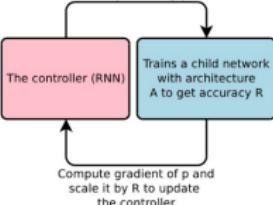
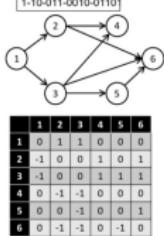
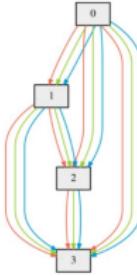
# Search Strategy

RL (Reinforcement Learning)	EA (Evolutionary Algorithms)																																										
																																											
(Sutton & Barto, 2018)	(Wistuba et al., 2019)																																										
 <p>Sample architecture A with probability <math>p</math></p> <p>The controller (RNN)</p> <p>Trains a child network with architecture A to get accuracy <math>R</math></p> <p>Compute gradient of <math>p</math> and scale it by <math>R</math> to update the controller</p>	 <p>1-10-011-0010-0110</p> <pre> graph LR     1((1)) --&gt; 2((2))     1((1)) --&gt; 3((3))     2((2)) --&gt; 4((4))     3((3)) --&gt; 4((4))     3((3)) --&gt; 5((5))     4((4)) --&gt; 6((6))     5((5)) --&gt; 6((6))     </pre> <table border="1"> <tr><td>1</td><td>2</td><td>3</td><td>4</td><td>5</td><td>6</td></tr> <tr><td>1</td><td>0</td><td>1</td><td>0</td><td>0</td><td>0</td></tr> <tr><td>2</td><td>-1</td><td>0</td><td>1</td><td>0</td><td>1</td></tr> <tr><td>3</td><td>-1</td><td>0</td><td>1</td><td>1</td><td>1</td></tr> <tr><td>4</td><td>0</td><td>-1</td><td>0</td><td>0</td><td>0</td></tr> <tr><td>5</td><td>0</td><td>0</td><td>-1</td><td>0</td><td>0</td></tr> <tr><td>6</td><td>0</td><td>-1</td><td>0</td><td>-1</td><td>0</td></tr> </table>	1	2	3	4	5	6	1	0	1	0	0	0	2	-1	0	1	0	1	3	-1	0	1	1	1	4	0	-1	0	0	0	5	0	0	-1	0	0	6	0	-1	0	-1	0
1	2	3	4	5	6																																						
1	0	1	0	0	0																																						
2	-1	0	1	0	1																																						
3	-1	0	1	1	1																																						
4	0	-1	0	0	0																																						
5	0	0	-1	0	0																																						
6	0	-1	0	-1	0																																						
(Zoph & Le, 2016)	(Lu et al., 2019)																																										

# Search Strategy

RL (Reinforcement Learning)	EA (Evolutionary Algorithms)	GB (Gradient-based)																																										
 <p>state <math>S_t</math> → Agent → action <math>A_t</math> → Environment → reward <math>R_t</math> → Agent</p> <p>(Sutton &amp; Barto, 2018)</p>	 <p>Initialization → Population → Parent Selection → Parents → Survivor Selection → Offspring → Recombination &amp; Mutation → Fitness Evaluation</p> <p>(Wistuba et al., 2019)</p>	 <p>"Locked" Network Parameters <math>\bar{\theta}</math> → Update Network <math>\nabla_{\bar{\theta}} L</math> → "Locked" Architecture Parameters <math>\bar{\alpha}</math> → Update Architecture <math>\nabla_{\bar{\alpha}} L</math></p> <p>(Green et al., 2019)</p>																																										
 <p>The controller (RNN) → Sample architecture A with probability <math>p</math> → Trains a child network with architecture A to get accuracy <math>R</math> → Compute gradient of <math>p</math> and scale it by <math>R</math> to update the controller</p> <p>(Zoph &amp; Le, 2016)</p>	 <p>DAG: 1 → 2, 1 → 3, 2 → 4, 3 → 4, 3 → 5, 4 → 6, 5 → 6.</p> <table border="1"> <tr><td>1</td><td>2</td><td>3</td><td>4</td><td>5</td><td>6</td></tr> <tr><td>1</td><td>0</td><td>1</td><td>0</td><td>0</td><td>0</td></tr> <tr><td>2</td><td>-1</td><td>0</td><td>1</td><td>0</td><td>1</td></tr> <tr><td>3</td><td>-1</td><td>0</td><td>1</td><td>1</td><td>1</td></tr> <tr><td>4</td><td>0</td><td>-1</td><td>0</td><td>0</td><td>0</td></tr> <tr><td>5</td><td>0</td><td>0</td><td>-1</td><td>0</td><td>0</td></tr> <tr><td>6</td><td>0</td><td>-1</td><td>0</td><td>0</td><td>-1</td></tr> </table> <p>(Lu et al., 2019)</p>	1	2	3	4	5	6	1	0	1	0	0	0	2	-1	0	1	0	1	3	-1	0	1	1	1	4	0	-1	0	0	0	5	0	0	-1	0	0	6	0	-1	0	0	-1	 <p>(Liu et al., 2019b)</p>
1	2	3	4	5	6																																							
1	0	1	0	0	0																																							
2	-1	0	1	0	1																																							
3	-1	0	1	1	1																																							
4	0	-1	0	0	0																																							
5	0	0	-1	0	0																																							
6	0	-1	0	0	-1																																							

# Search Strategy

RL (Reinforcement Learning)	EA (Evolutionary Algorithms)	GB (Gradient-based)	Others																																										
 <p>state <math>S_t</math></p> <p>reward <math>R_t</math></p> <p>action <math>A_t</math></p> <p>Environment</p> <p>Agent</p>	 <p>Initialization</p> <p>Population</p> <p>Parent Selection</p> <p>Survivor Selection</p> <p>Recombination &amp; Mutation</p> <p>Offspring</p> <p>Fitness Evaluation</p> <p>Update Network <math>\nabla_{\theta} L</math></p>	 <p>"Locked" Network Parameters <math>\theta</math></p> <p>"Locked" Architecture Parameters <math>\alpha</math></p> <p>Update Network <math>\nabla_{\theta} L</math></p> <p>Update Architecture <math>\nabla_{\alpha} L</math></p>	<ul style="list-style-type: none"> <li>• Bayesian Optimization</li> <li>• Random Search</li> <li>• Model-based <ul style="list-style-type: none"> <li>◦ Gaussian process</li> <li>◦ Random Forest</li> <li>◦ XGBoost</li> </ul> </li> </ul>																																										
(Sutton & Barto, 2018)	(Wistuba et al., 2019)	(Green et al., 2019)																																											
 <p>Sample architecture A with probability <math>p</math></p> <p>The controller (RNN)</p> <p>Trains a child network with architecture A to get accuracy <math>R</math></p> <p>Compute gradient of <math>p</math> and scale it by <math>R</math> to update the controller</p>	 <p>1-10-011-0010-0110</p> <p>1 2 3 4 5 6</p> <table border="1"> <tr><td>1</td><td>0</td><td>1</td><td>1</td><td>0</td><td>0</td><td>0</td></tr> <tr><td>2</td><td>-1</td><td>0</td><td>1</td><td>0</td><td>1</td><td></td></tr> <tr><td>3</td><td>-1</td><td>0</td><td>0</td><td>1</td><td>1</td><td>1</td></tr> <tr><td>4</td><td>0</td><td>-1</td><td>-1</td><td>0</td><td>0</td><td>0</td></tr> <tr><td>5</td><td>0</td><td>0</td><td>-1</td><td>0</td><td>0</td><td>1</td></tr> <tr><td>6</td><td>0</td><td>0</td><td>-1</td><td>0</td><td>-1</td><td>0</td></tr> </table>	1	0	1	1	0	0	0	2	-1	0	1	0	1		3	-1	0	0	1	1	1	4	0	-1	-1	0	0	0	5	0	0	-1	0	0	1	6	0	0	-1	0	-1	0	 <p>0</p> <p>1</p> <p>2</p> <p>3</p>	<ul style="list-style-type: none"> <li>(Zoph &amp; Le, 2016)</li> <li>(Lu et al., 2019)</li> <li>(Liu et al., 2019b)</li> <li>(Luo et al., 2019)</li> </ul>
1	0	1	1	0	0	0																																							
2	-1	0	1	0	1																																								
3	-1	0	0	1	1	1																																							
4	0	-1	-1	0	0	0																																							
5	0	0	-1	0	0	1																																							
6	0	0	-1	0	-1	0																																							

# NAS framework

- **Search Space:**

- Cell-based
- One-Shot

- **Search Strategy:**

- Reinforcement Learning
- Gradient-based
- Evolutionary Algorithms

- **Performance Estimation:**

- Full training
- Partial training
- Performance prediction

# Once-for-All

- **Search Space:**

- Cell-based
- One-Shot → Training stage

- **Search Strategy:**

- Reinforcement Learning
- Gradient-based → Training stage
- Evolutionary Algorithms → Search stage

- **Performance Estimation:**

- Full training
- Partial training
- Performance prediction → Search stage

# Agenda

- 1 Introduction
- 2 AutoML
- 3 Training stage
- 4 Search stage
  - Random search
  - OFA search
  - OFA<sup>2</sup> search
  - Comparison
- 5 Ensemble
  - Manual sampling
  - Multi-objective Optimization
- 6 Conclusion
- 7 References

# OFA network

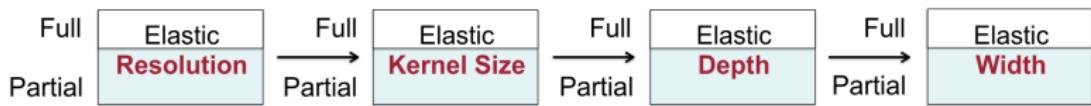
- Convolutional Neural Network (CNN)
  - Same architecture space of MobileNetV3
  - Supernet with smaller subnetworks nested inside larger architectures
- 
- Final architecture: 5 units
  - Each unit have 3 adjustable parameters
    - depth  $\in \{2, 3, 4\}$
    - width  $\in \{3, 4, 6\}$
    - kernel size  $\in \{3, 5, 7\}$
  - The 4<sub>th</sub> parameter is unique for the architecture
    - image resolution  $\in \{128, 132, \dots, 224\}$

# OFA training

- Start training the network at its full size
  - Optimizer = SGD
  - Nesterov momentum = 0.9
  - Weight decay =  $3e^{-5}$
  - Initial learning rate = 2.6
  - Learning weight decay = cosine schedule
  - Epochs = 180
  - Batch size = 2048
  - GPUs = 32 V100 GPUs
  - Cost = 1,200 GPU hours

(Cai et al. 2020).

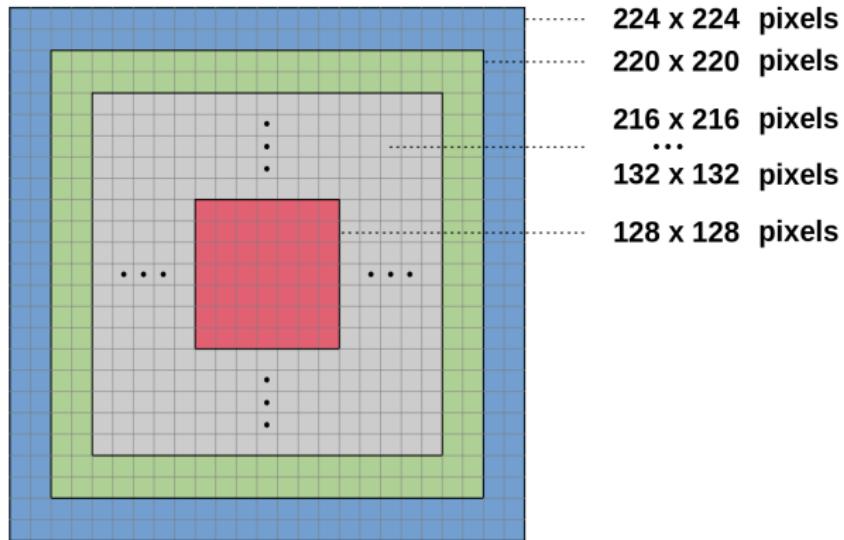
# Progressive shrinking



- Progressively shrinks the network using the *Progressive shrinking* algorithm

# Elastic resolution

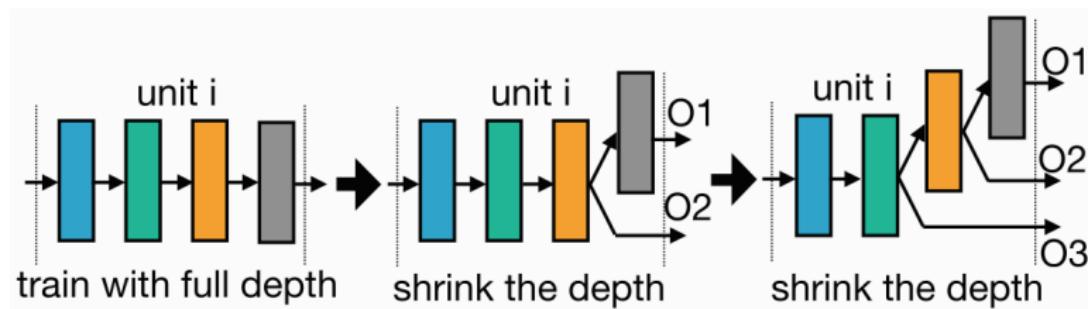
- 25 different input resolutions:  $\{128, 132, \dots, 224\}$



(Cai et al. 2020)

## Elastic depth

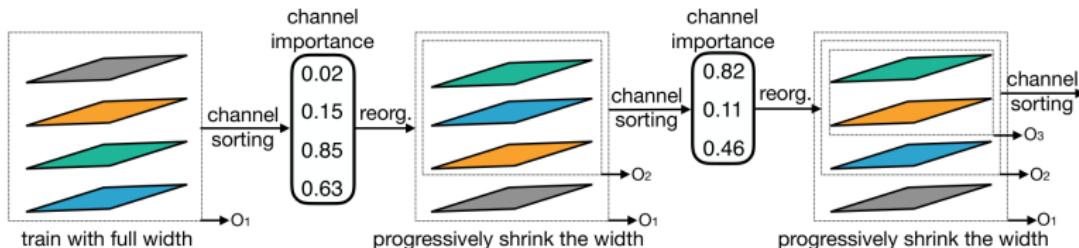
- Depth options: {2, 3, 4} layers.
- Start training with 4 layers, then trains with 3 layer, and lastly with 2 layers.



(Cai et al. 2020)

# Elastic width

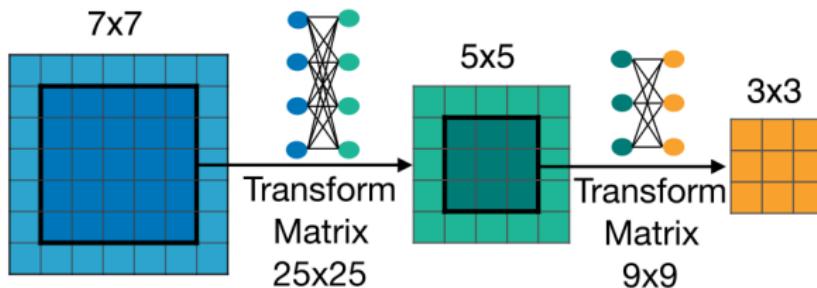
- Channel expansion ratio: {3, 4, 6}.
- Start training with all filters. Then skip the least important one by one.
- The importance of each channel is calculated using the L1 norm of the channel's weights.



(Cai et al. 2020)

## Elastic kernel size

- Convolutional kernel size:  $\{7 \times 7, 5 \times 5, 3 \times 3\}$
- Start training with kernel  $7 \times 7$ , then  $5 \times 5$ , and lastly  $3 \times 3$ .



(Cai et al. 2020)

- Different layers have different kernel transformation matrices.
- The matrices are shared among different channels within each layer.
- Total kernel transformation parameters:  $25 \times 25 + 9 \times 9 = 706$  per layer

# Agenda

- 1 Introduction
- 2 AutoML
- 3 Training stage
- 4 Search stage
  - Random search
  - OFA search
  - OFA<sup>2</sup> search
  - Comparison
- 5 Ensemble
  - Manual sampling
  - Multi-objective Optimization
- 6 Conclusion
- 7 References

# Performance estimator

- **Accuracy Predictor**

- Feedforward neural network with
  - 3 layers
  - 400 hidden units per layer

- **Efficiency Predictor**

- Latency lookup table (Cai, Zhu, and Han 2022)
- Function to count FLOPS given the operations.

# Architecture Space

PS	property	options	# choices
elastic resolution	input resolution	{224, ..., 128}	25
elastic kernel	kernel size	{3, 5, 7}	3
elastic depth	layers	{2, 3, 4}	3
elastic width	channels	{3, 4, 6}	3

# Architecture Space

PS	property	options	# choices
elastic resolution	input resolution	{224, ..., 128}	25
elastic kernel	kernel size	{3, 5, 7}	3
elastic depth	layers	{2, 3, 4}	3
elastic width	channels	{3, 4, 6}	3

```
{'ks': [5, 7, 3, 5, 5, 7, 7, 5, 3, 3, 3, 3, 7, 5, 3, 7, 7, 5, 7, 5],
 'e': [4, 4, 4, 3, 3, 4, 3, 4, 6, 6, 4, 3, 6, 3, 6, 4, 3, 3, 4, 3],
 'd': [2, 4, 2, 3, 2],
 'r': [176]}
```

# Architecture Space

PS	property	options	# choices
elastic resolution	input resolution	{224, ..., 128}	25
elastic kernel	kernel size	{3, 5, 7}	3
elastic depth	layers	{2, 3, 4}	3
elastic width	channels	{3, 4, 6}	3

```

{'ks': [5, 7, 3, 5, 5, 7, 7, 5, 3, 3, 3, 3, 7, 5, 3, 7, 7, 5, 7, 5],
 'e': [4, 4, 4, 3, 3, 4, 3, 4, 6, 6, 4, 3, 6, 3, 6, 4, 3, 3, 4, 3],
 'd': [2, 4, 2, 3, 2],
 'r': [176]}

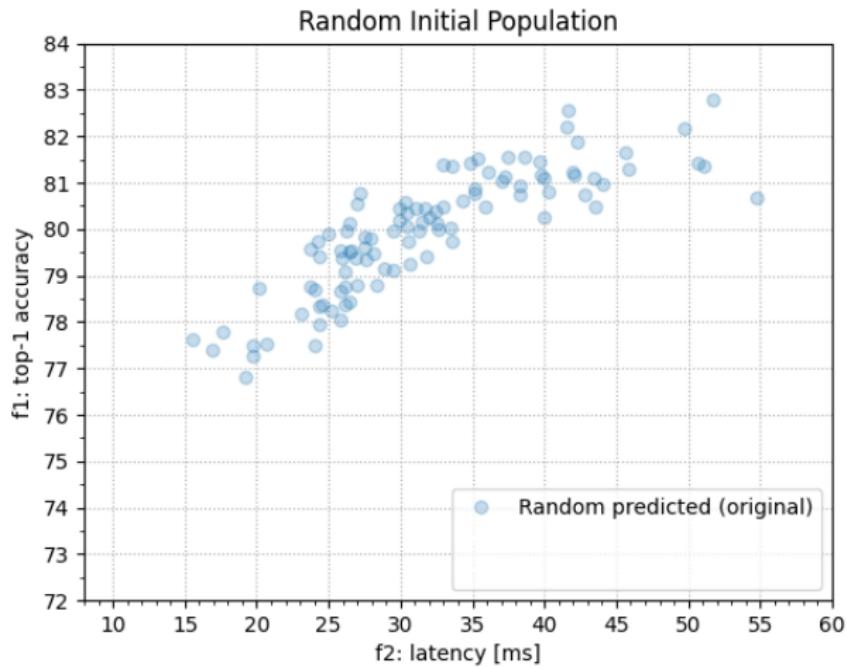
```

- Search space of  $((3 \times 3)^2 + (3 \times 3)^3 + (3 \times 3)^4)^5 \approx 2 \times 10^{19}$  different neural network architectures.
- All these subnetworks share the same weights: total of 7.7M parameters.

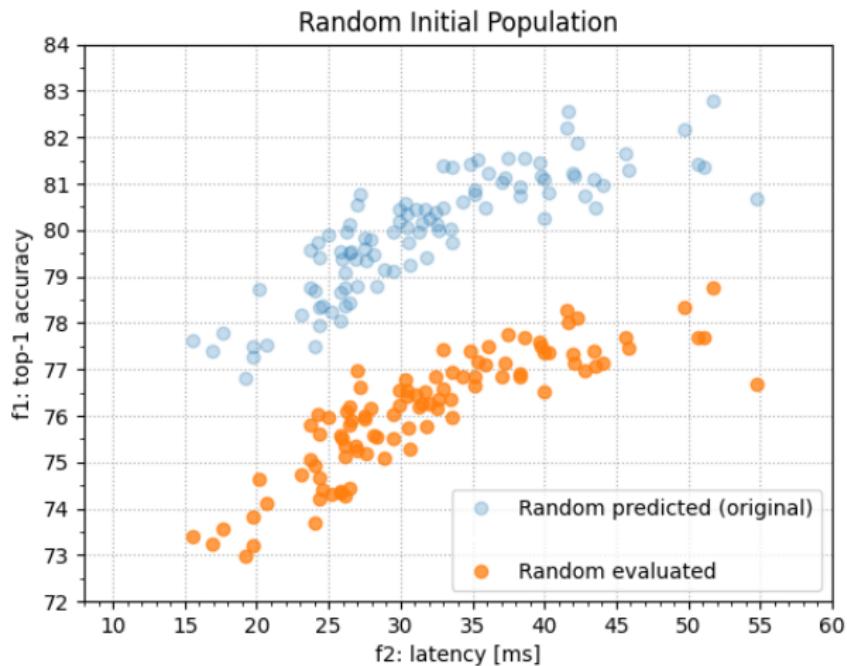
# Agenda

- 1 Introduction
- 2 AutoML
- 3 Training stage
- 4 Search stage
  - Random search
  - OFA search
  - OFA<sup>2</sup> search
  - Comparison
- 5 Ensemble
  - Manual sampling
  - Multi-objective Optimization
- 6 Conclusion
- 7 References

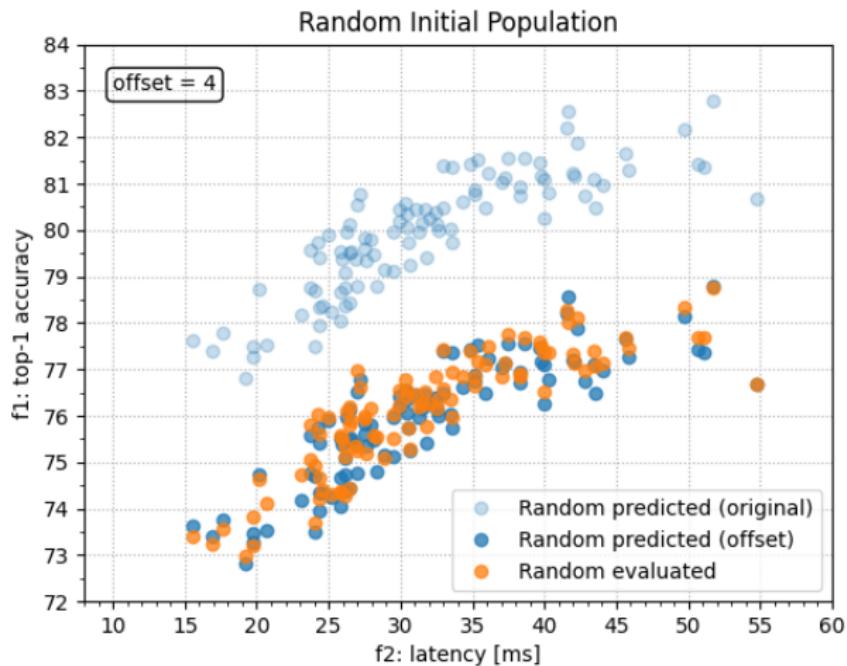
# Random search



## Random search



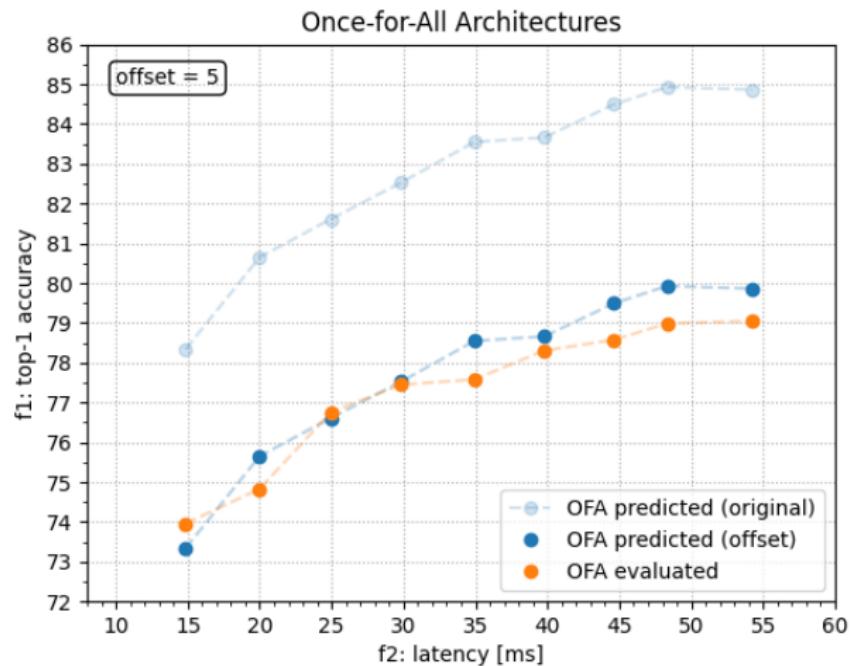
## Random search



# Agenda

- 1 Introduction
- 2 AutoML
- 3 Training stage
- 4 Search stage
  - Random search
  - **OFA search**
  - OFA<sup>2</sup> search
  - Comparison
- 5 Ensemble
  - Manual sampling
  - Multi-objective Optimization
- 6 Conclusion
- 7 References

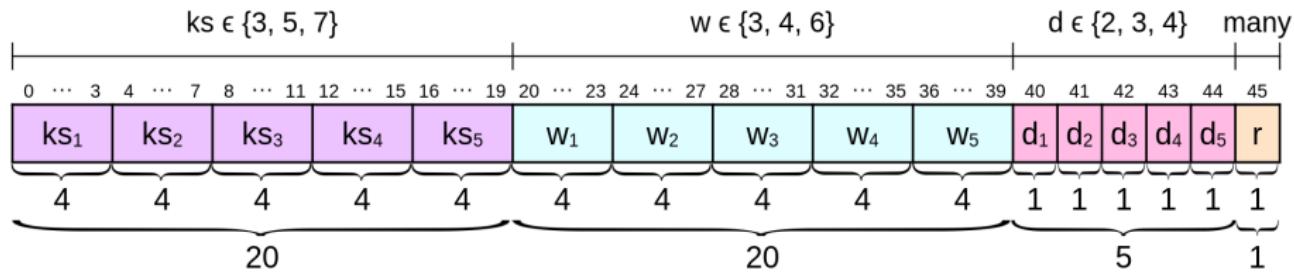
## OFA search



# Agenda

- 1 Introduction
- 2 AutoML
- 3 Training stage
- 4 Search stage
  - Random search
  - OFA search
  - OFA<sup>2</sup> search
  - Comparison
- 5 Ensemble
  - Manual sampling
  - Multi-objective Optimization
- 6 Conclusion
- 7 References

# Encoding



# Encoding example

0	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18	19	20	21	22	23	24	25	26	27	28	29	30	31	32	33	34	35	36	37	38	39	40	41	42	43	44	45
3	3	7	5	3	7	3	5	3	7	5	3	7	3	5	5	3	6	3	4	4	3	4	3	6	4	4	3	3	4	3	3	4	6	3	3	3	4	2	2	3	224				

# Encoding example

0 1 2 3 4 5 6 7 8 9 10 11 12 13 14 15 16 17 18 19 20 21 22 23 24 25 26 27 28 29 30 31 32 33 34 35 36 37 38 39 40 41 42 43 44 45

3	3	7	5	7	3	5	3	3	3	7	5	7	5	3	7	3	5	5	3	4	2	2	3	224
---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	-----

													size	index								
ks	3	3	7	5	7	3	5	3	3	3	7	5	3	7	3	5	5	3	20	[00:19]		
w	6	3	4	4	3	4	3	6	4	4	3	3	4	3	3	3	4	6	3	3	20	[20:39]
d	3				4				2			2			2		3		5	[40:44]		
r	224																	1		[45]		

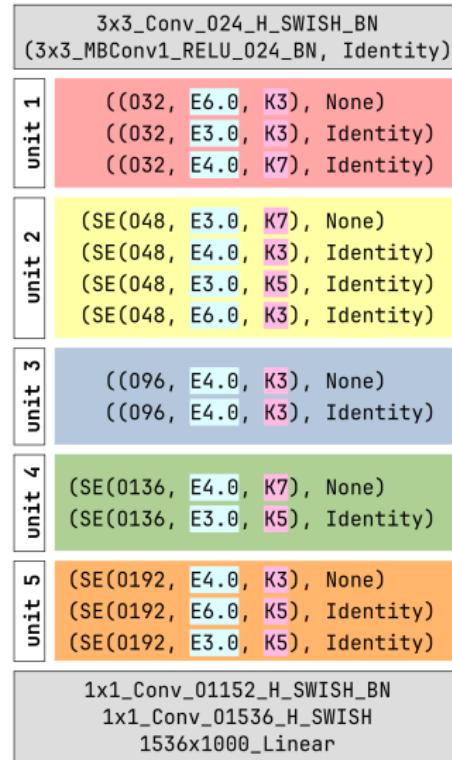
## Encoding example

0	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18	19	20	21	22	23	24	25	26	27	28	29	30	31	32	33	34	35	36	37	38	39	40	41	42	43	44	45
3375	7353	3375	7537	3553	6344	3436	4433	4333	4633	3	4	2	2	3	224																														

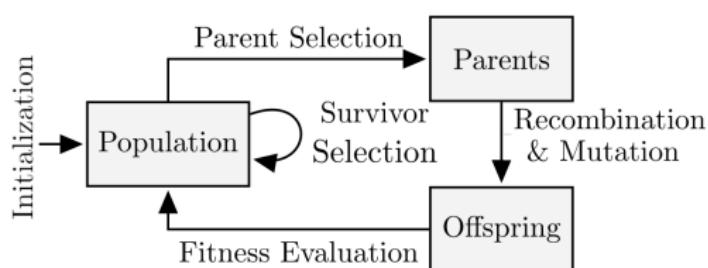
	size											index											
ks	3	3	7	5	7	3	5	3	3	3	7	5	7	5	3	7	3	5	5	3	20	[00:19]	
w	6	3	4	4	3	4	3	6	4	4	3	3	4	3	3	3	4	6	3	3	20	[20:39]	
d	3		4			2			2			3			5	[40:44]							
r	224																					1	[45]

	Unit 1	Unit 2	Unit 3	Unit 4	Unit 5	size	index
ks	3 3 7 	7 3 5 3	3 3  	7 5  	3 5 5 	20	[00:19]
w	6 3 4 	3 4 3 6	4 4  	4 3  	4 6 3 	20	[20:39]
d	3 	4	2  	2  	3 	5	[40:44]
r	224					1	[45]

# Encoding example



# Evolutionary Algorithm framework



(Real et al. 2017)

## Elements:

- **Individual:** candidate solution
- **Encoding:** representation of a candidate solution
- **Fitness function:** evaluation of the candidate solution
- **Population:** group of individuals
- **Operators:**
  - Sampling
  - Mutation
  - Crossover
  - Selection
- **Offspring:** next generation population

# Operator: Mutation

Original individual

0	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18	19	20	21	22	23	24	25	26	27	28	29	30	31	32	33	34	35	36	37	38	39	40	41	42	43	44	45
3	3	7	5	3	3	3	7	5	3	7	5	3	7	5	3	3	5	5	3	6	3	4	4	3	4	3	6	4	4	3	3	3	4	6	3	3	3	3	4	2	2	3	2	24	

The diagram illustrates the mutation process. Four arrows point downwards from the original individual's genome to the mutated individual's genome, specifically targeting the 10th, 21st, 37th, and 44th positions. These positions are highlighted in red in both the original and mutated genomes.

0	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18	19	20	21	22	23	24	25	26	27	28	29	30	31	32	33	34	35	36	37	38	39	40	41	42	43	44	45
3	3	7	5	3	3	3	7	5	3	7	5	3	7	5	3	3	5	5	3	6	3	4	4	3	4	3	6	4	4	3	3	3	4	6	3	3	3	3	4	2	2	3	2	24	

After mutation

# Operator: Crossover (recombination)

Parent 1

0 1 2 3 4 5 6 7 8 9 10 11 12 13 14 15 16 17 18 19 20 21 22 23 24 25 26 27 28 29 30 31 32 33 34 35 36 37 38 39 40 41 42 43 44 45

3	3	7	5	7	3	5	3	3	7	5	3	7	5	3	6	3	4	4	3	4	3	4	3	3	3	4	6	3	3	3	4	6	3	3	3	4	2	2	3	2	24
---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	----

3	7	5	5	7	5	3	3	3	5	7	5	3	5	3	7	3	3	6	3	4	4	3	6	3	3	3	4	3	3	3	4	6	3	3	2	4	2	4	3	2	24
---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	----

7	7	5	3	3	5	7	3	5	3	7	3	3	6	3	4	6	4	3	6	3	3	3	3	4	6	6	3	4	4	3	6	2	3	3	4	3	2	24
---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	----

0 1 2 3 4 5 6 7 8 9 10 11 12 13 14 15 16 17 18 19 20 21 22 23 24 25 26 27 28 29 30 31 32 33 34 35 36 37 38 39 40 41 42 43 44 45

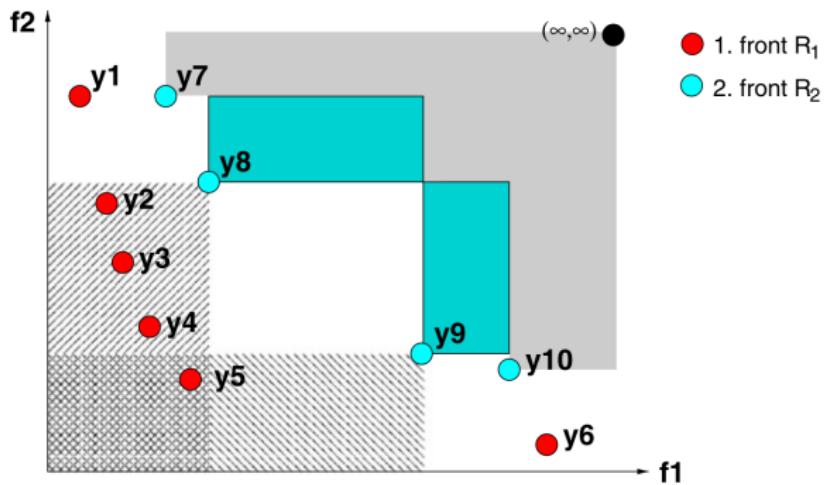
Parent 2

# Hyperparameters

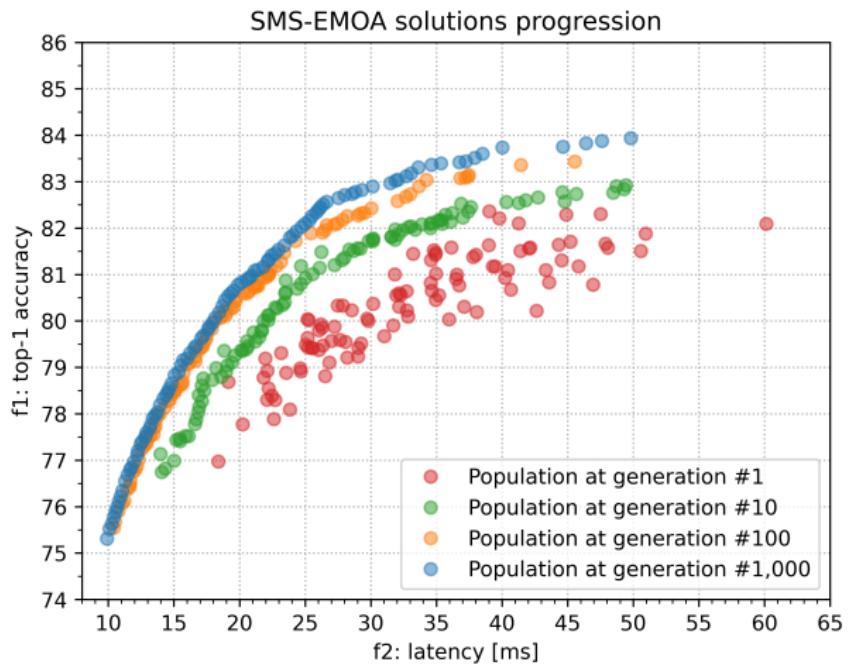
- Population: 100
- Generations: 1,000
- Operators:
  - Mutation: 50% of population
    - Mutation architecture (bitflip): 10%
  - Crossover: 25% as parents
    - Uniform
- Algorithms:
  - NSGA-II
  - SMS-EMOA
- Fitness function:
  - Accuracy predictor
  - Efficiency predictor
    - Latency
    - FLOPS

## SMS-EMOA

- Maximizes the hypervolume

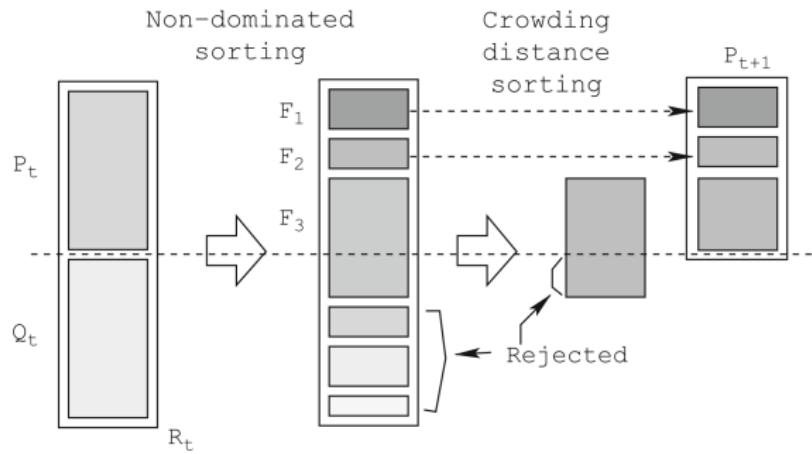


# SMS-EMOA: Population

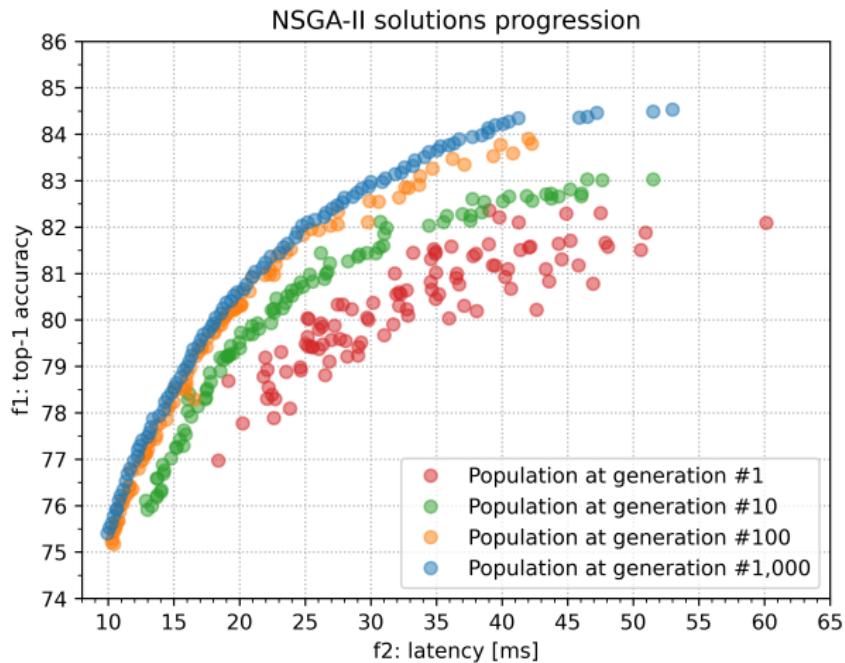


# NSGA-II

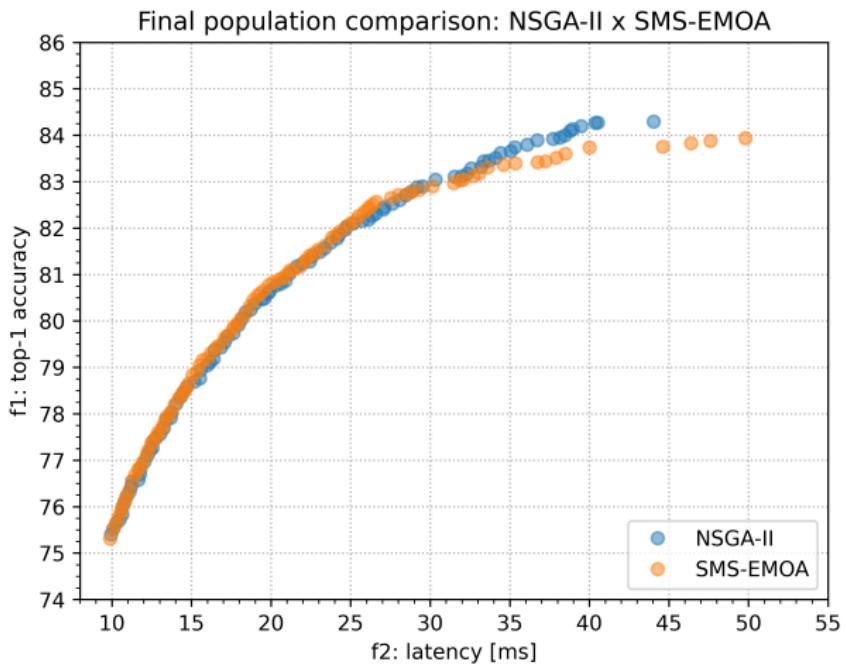
- **Sorting:**
  - Non-dominated
  - Crowding distance



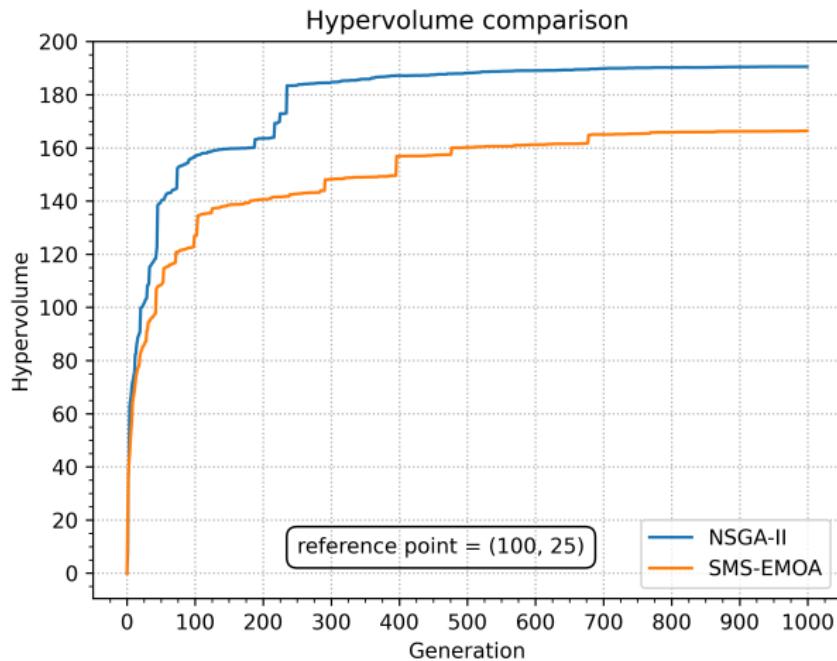
## NSGA-II: Population



# Population comparison



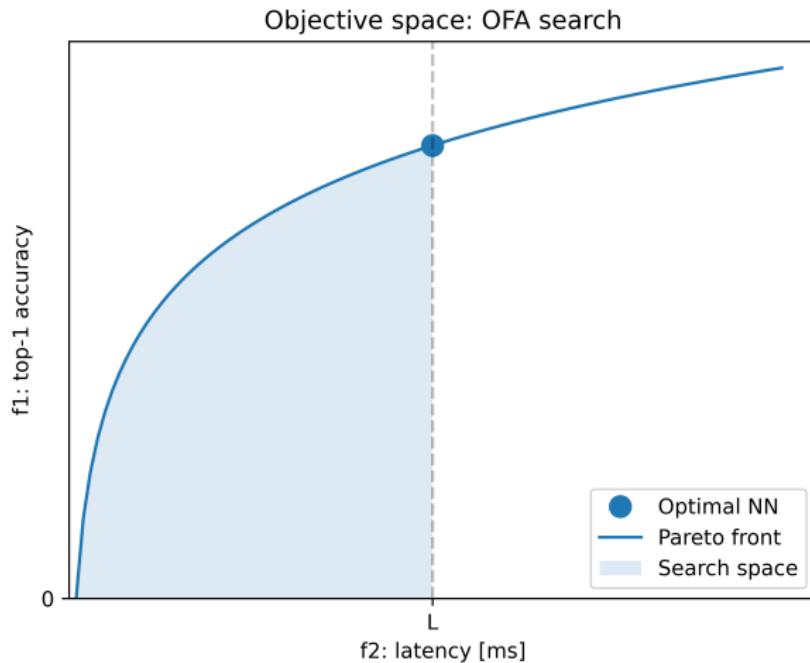
# Hypervolume



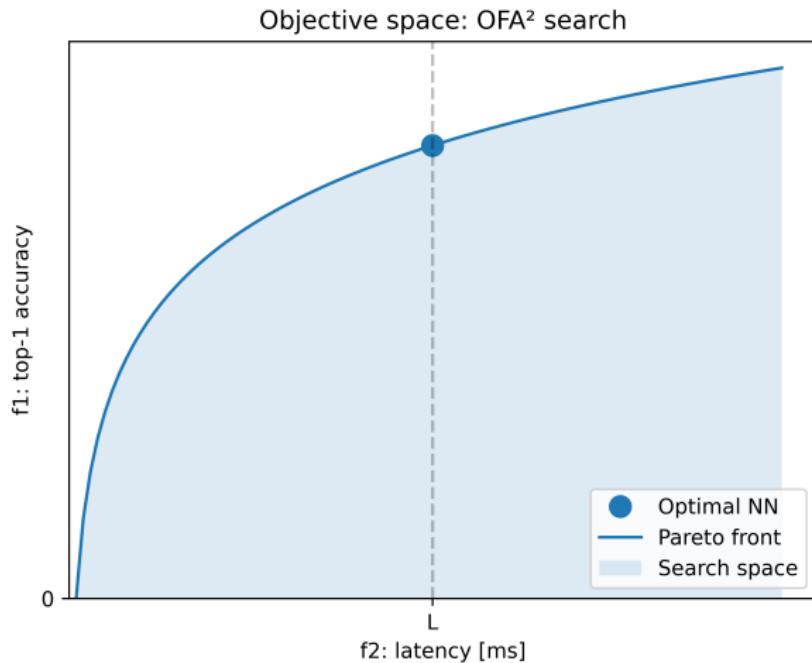
# Agenda

- 1 Introduction
- 2 AutoML
- 3 Training stage
- 4 Search stage
  - Random search
  - OFA search
  - OFA<sup>2</sup> search
  - Comparison
- 5 Ensemble
  - Manual sampling
  - Multi-objective Optimization
- 6 Conclusion
- 7 References

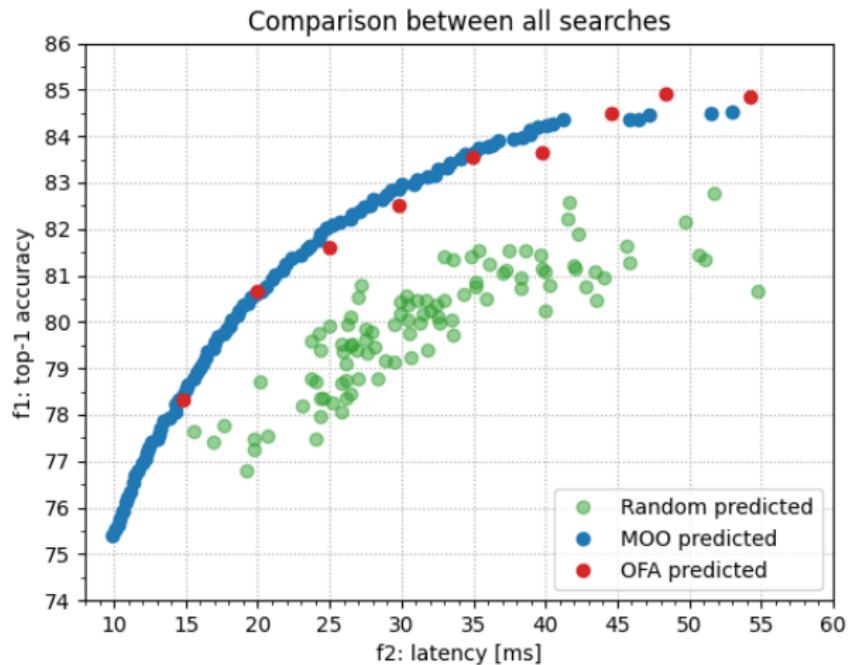
## OFA search space



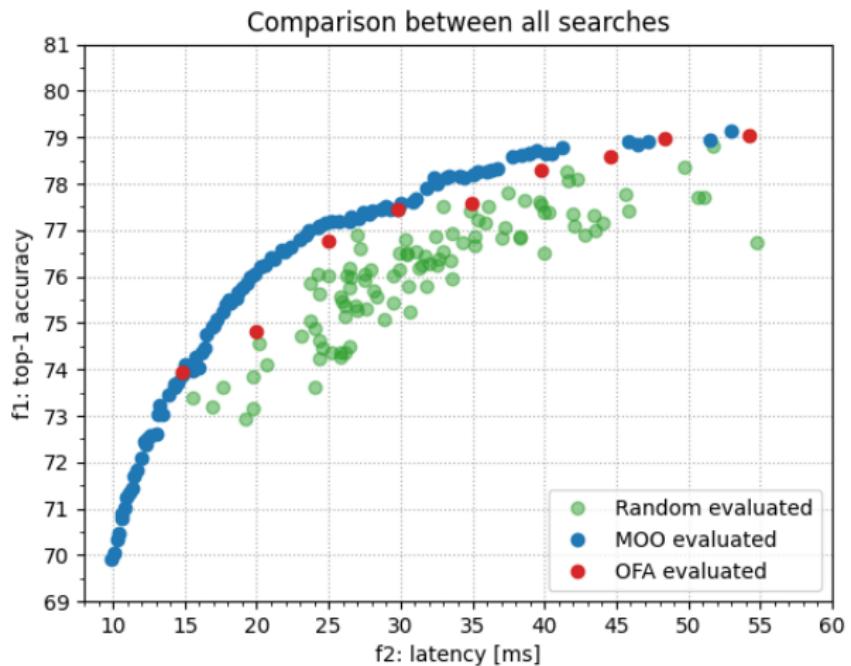
# OFA<sup>2</sup> search space



# Predicted accuracy



## Evaluated accuracy



# Agenda

- 1 Introduction
- 2 AutoML
- 3 Training stage
- 4 Search stage
  - Random search
  - OFA search
  - OFA<sup>2</sup> search
  - Comparison
- 5 Ensemble
  - Manual sampling
  - Multi-objective Optimization
- 6 Conclusion
- 7 References

# Overview

- **Voting:**
  - Hard voting
  - Soft voting
- **Latency:**
  - Sum
  - Maximum
- **Formation:**
  - Manual Sampling
  - Multi-Objective Optimization

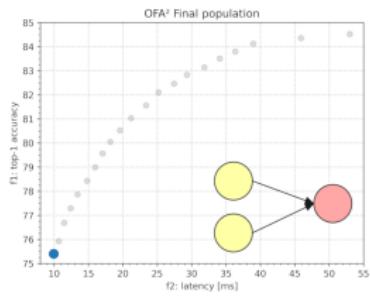
# Voting

- **Voting:**

- Hard voting
- Soft voting

# Hard voting

Component 1



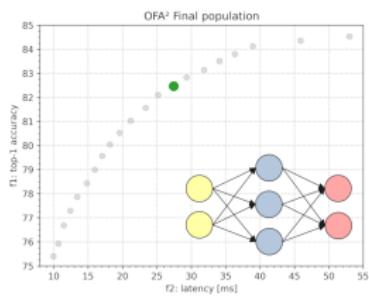
arg max

Class 7



Input image

Component 2

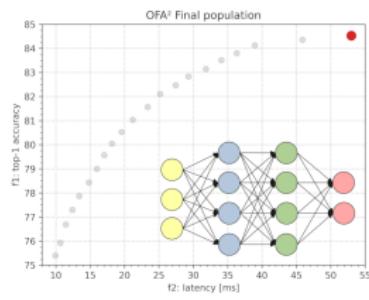


arg max

Class 7

Ensemble output: Class 7

Component 3

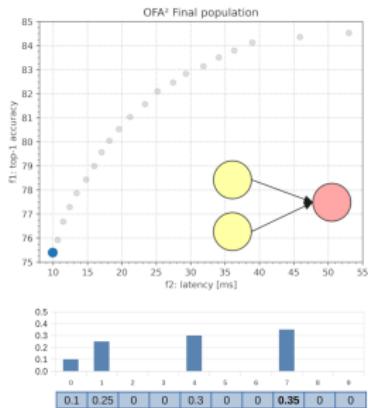


arg max

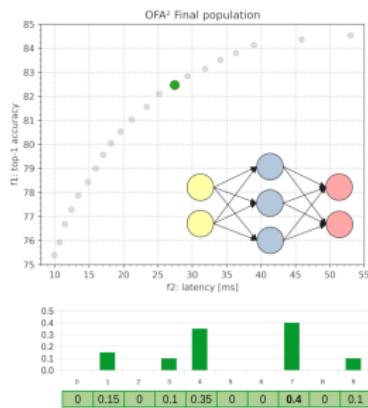
Class 4

# Soft voting

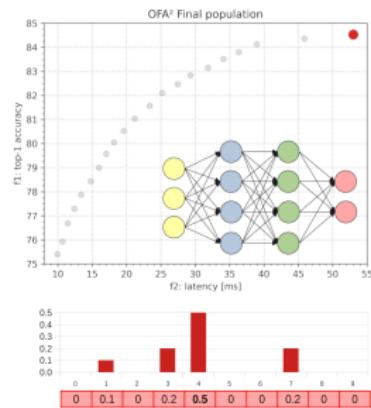
Component 1



Component 2

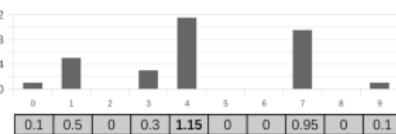


Component 3



Input image

Ensemble output:



$\xrightarrow{\text{arg max}}$  Class 4

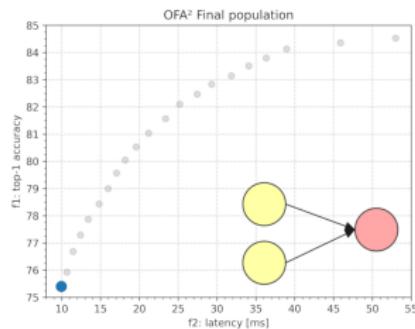
# Latency

- **Latency:**

- Sum
- Maximum

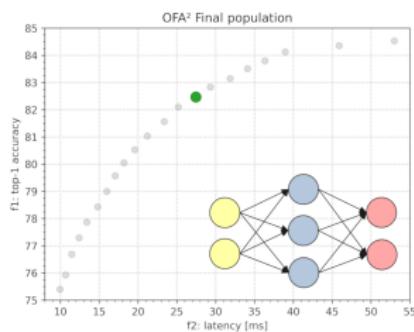
## Sum

Component 1



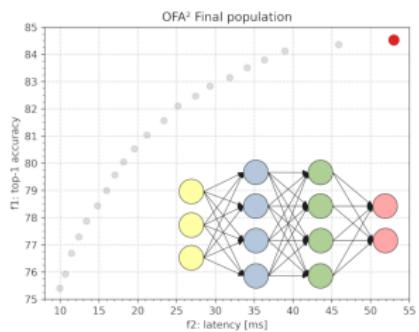
Latency = 10 ms

Component 2



Latency = 30 ms

Component 3



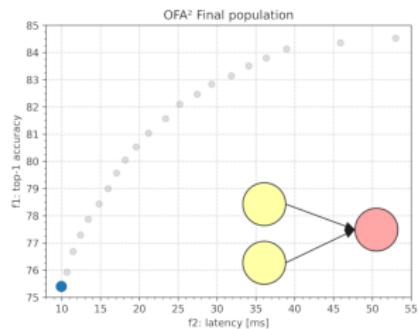
Latency = 50 ms



**Ensemble latency = 90 ms**

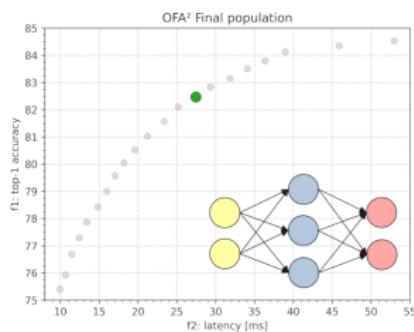
# Maximum

Component 1



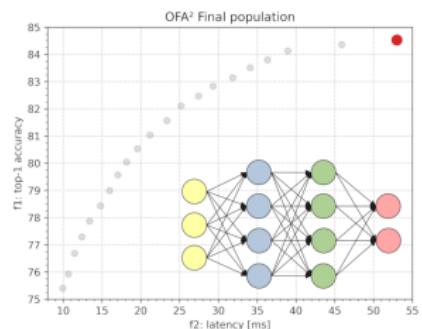
Latency = 10 ms

Component 2



Latency = 30 ms

Component 3



Latency = 50 ms

MAX

Ensemble latency = 50 ms

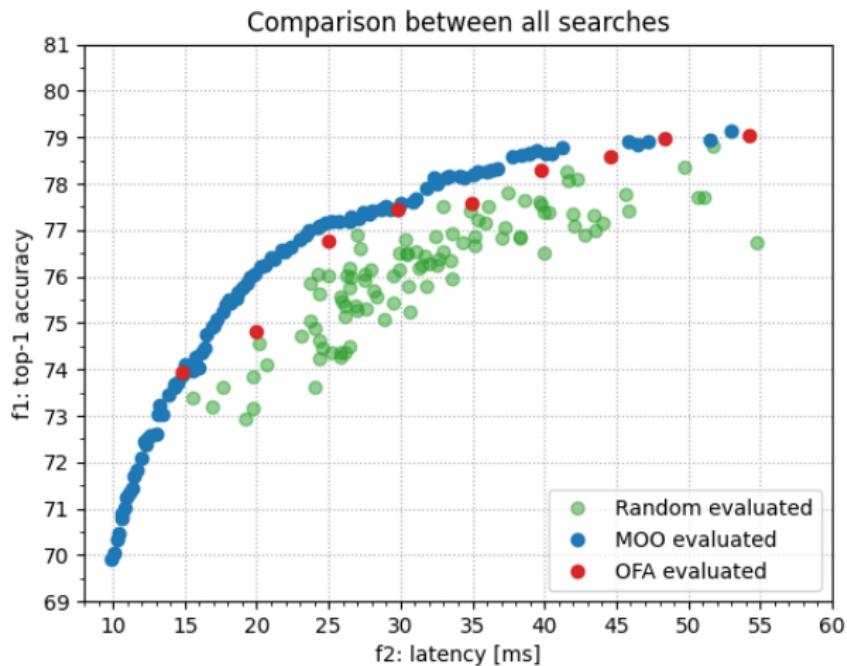
# Agenda

- 1 Introduction
- 2 AutoML
- 3 Training stage
- 4 Search stage
  - Random search
  - OFA search
  - OFA<sup>2</sup> search
  - Comparison
- 5 Ensemble
  - **Manual sampling**
  - Multi-objective Optimization
- 6 Conclusion
- 7 References

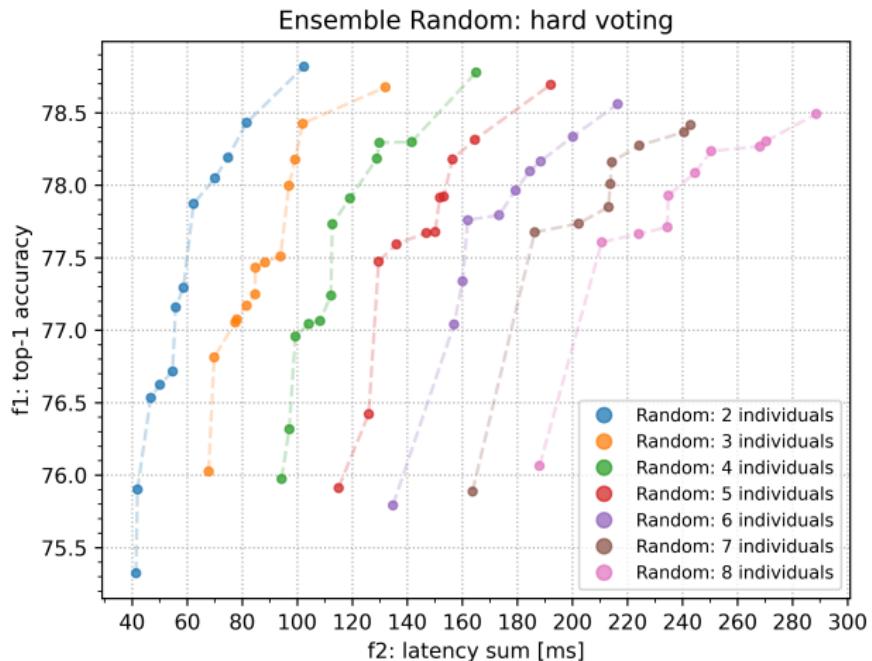
# Overview

- **Experiments:**
  - 43 ensembles for each scenario
- **Individuals:**
  - 2 individuals
  - 3 individuals
  - 4 individuals
  - 5 individuals
  - 6 individuals
  - 7 individuals
  - 8 individuals
- **Voting:**
  - Hard voting
  - Soft voting
- **Latency:**
  - Sum
  - Maximum

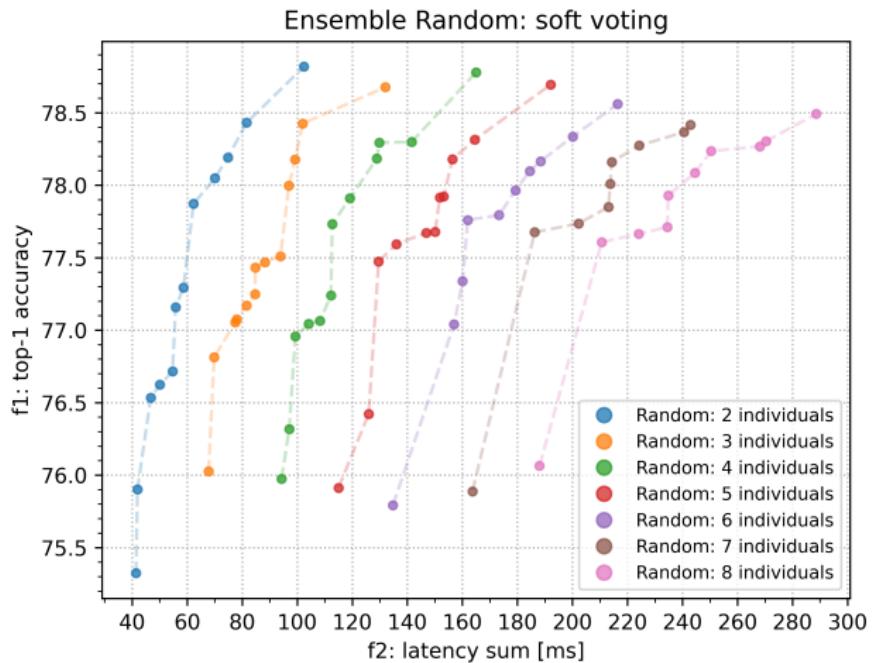
# Random search



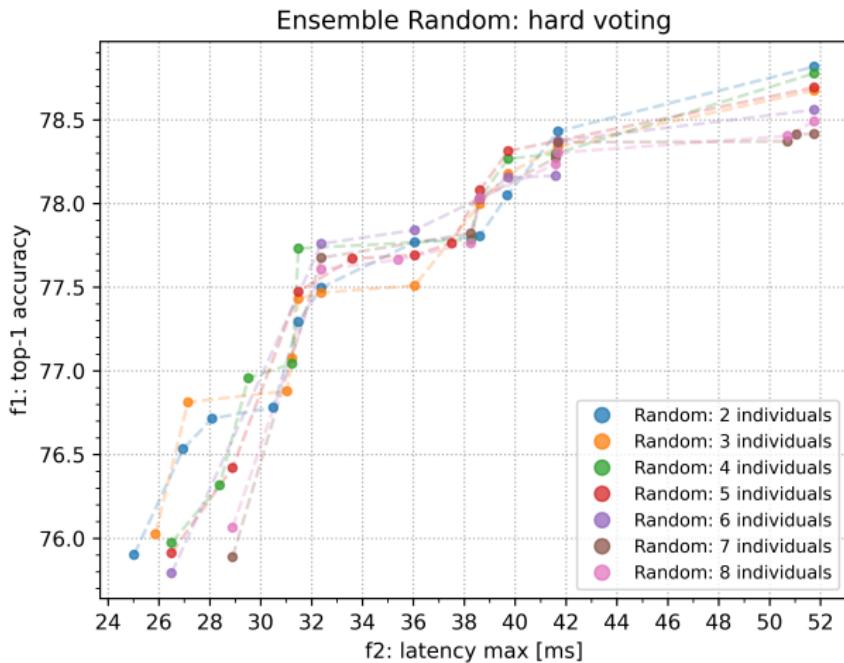
## Random search: hard, sum



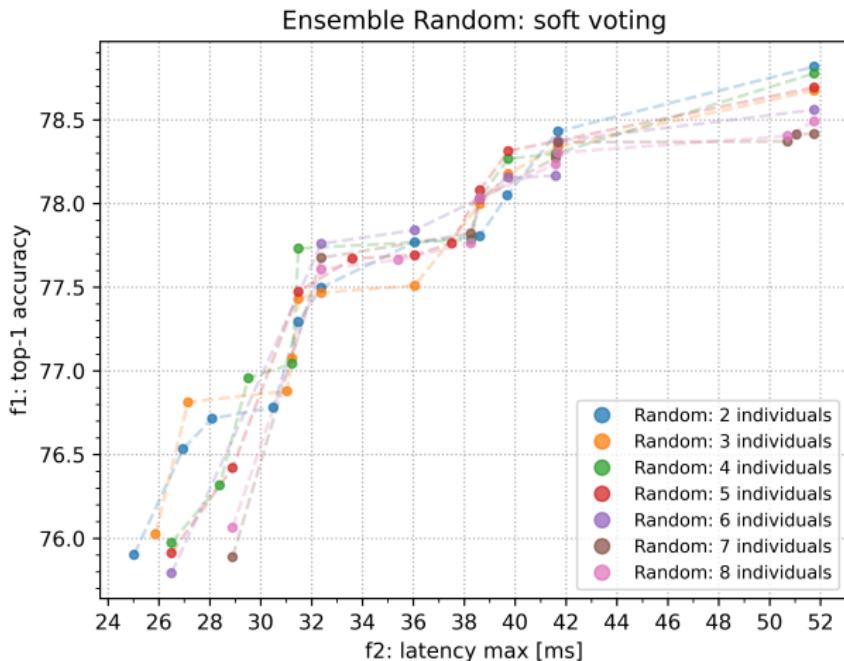
## Random search: soft, sum



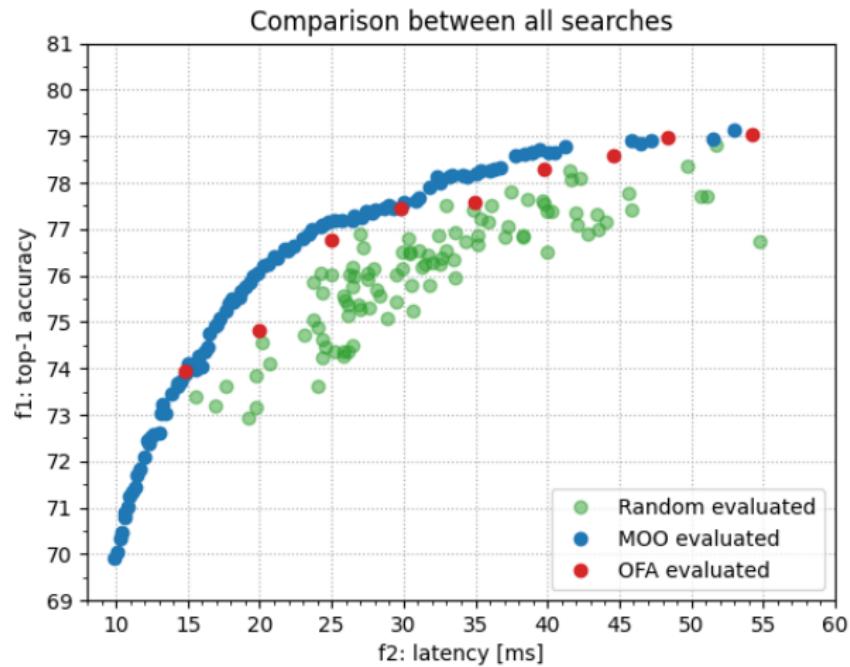
## Random search: hard, max



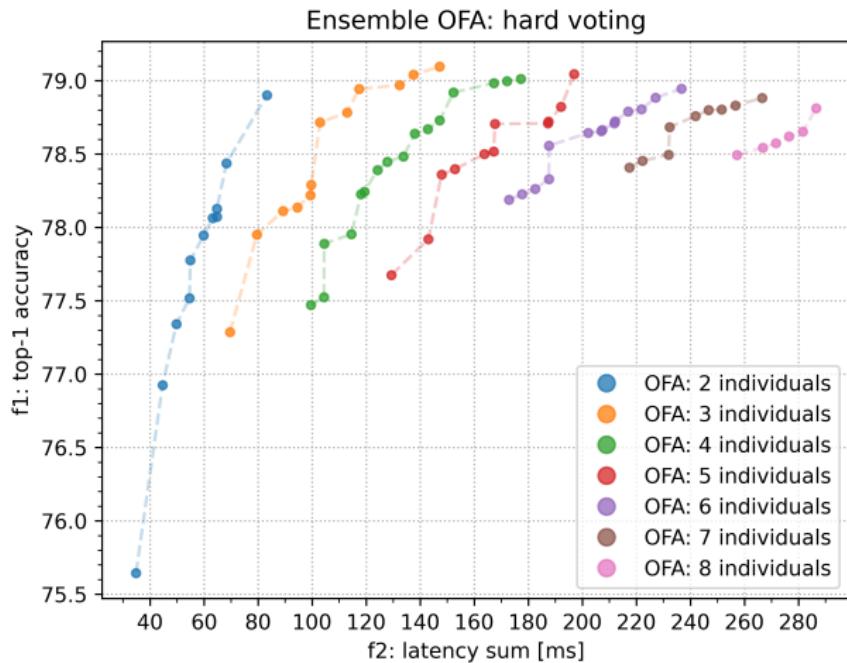
## Random search: soft, max



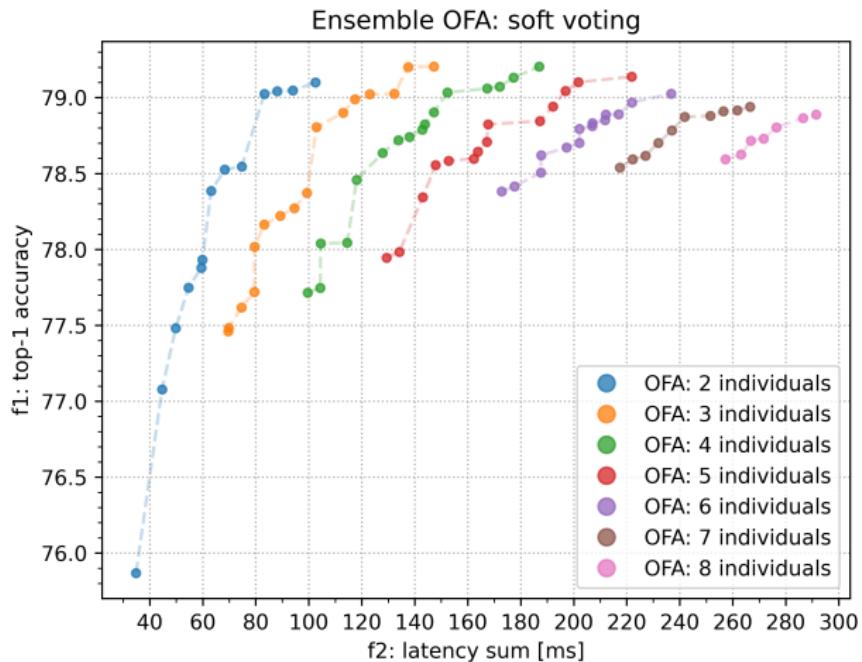
## OFA search



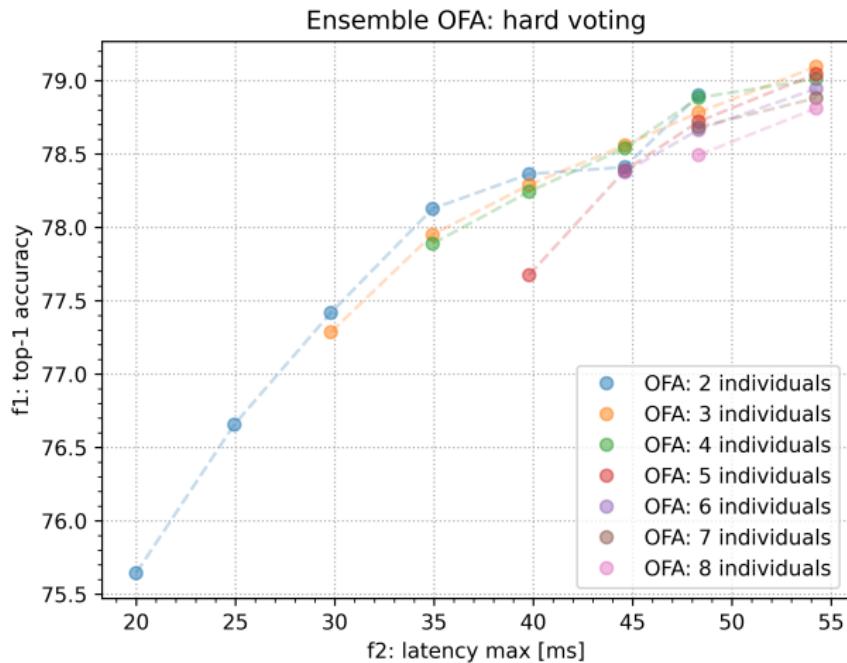
## OFA search: hard, sum



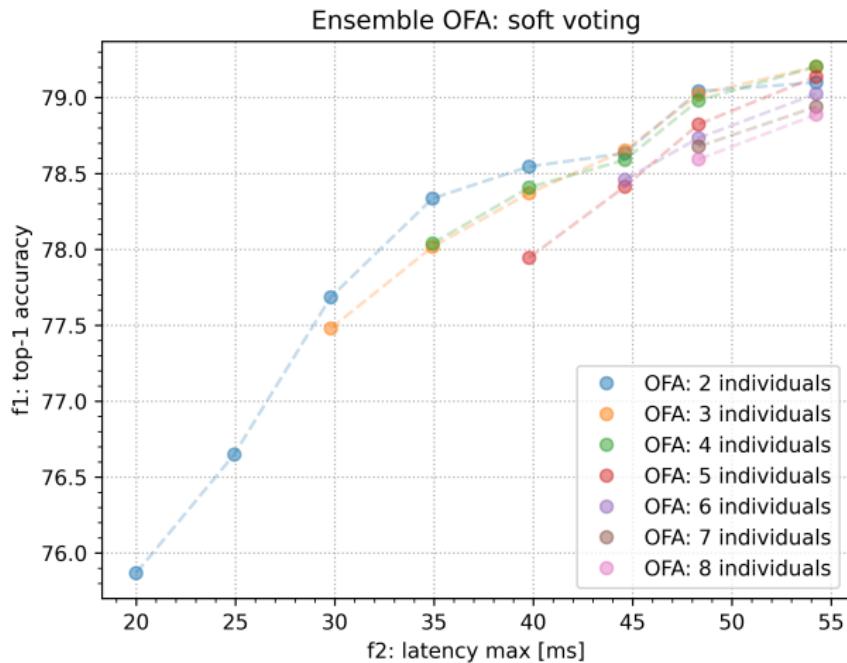
## OFA search: soft, sum

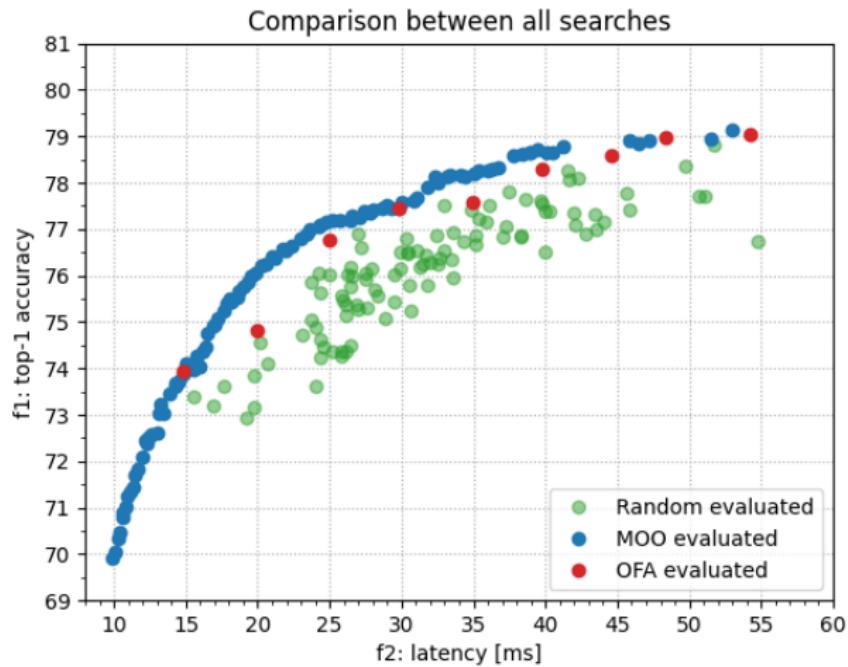


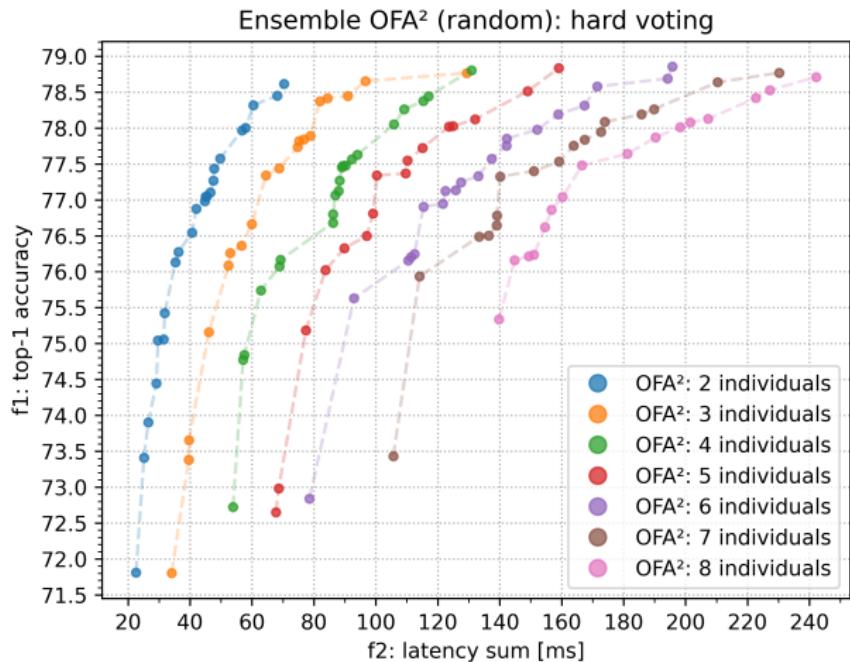
## OFA search: hard, max

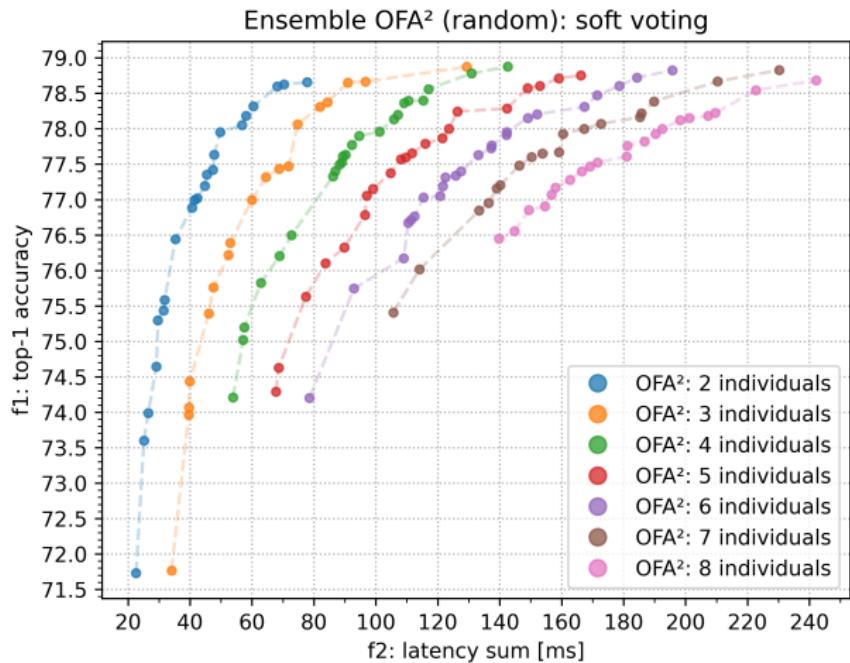


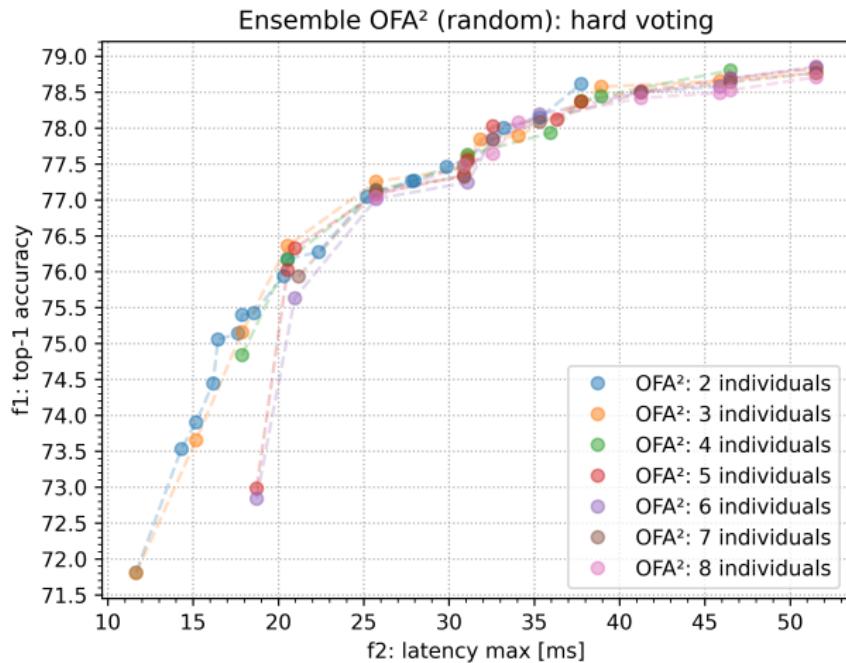
## OFA search: soft, max

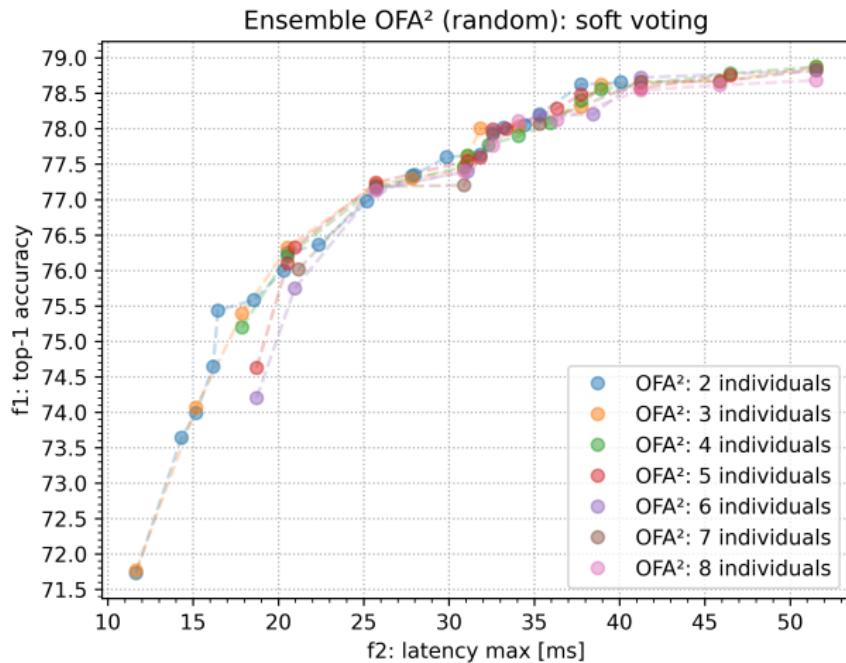


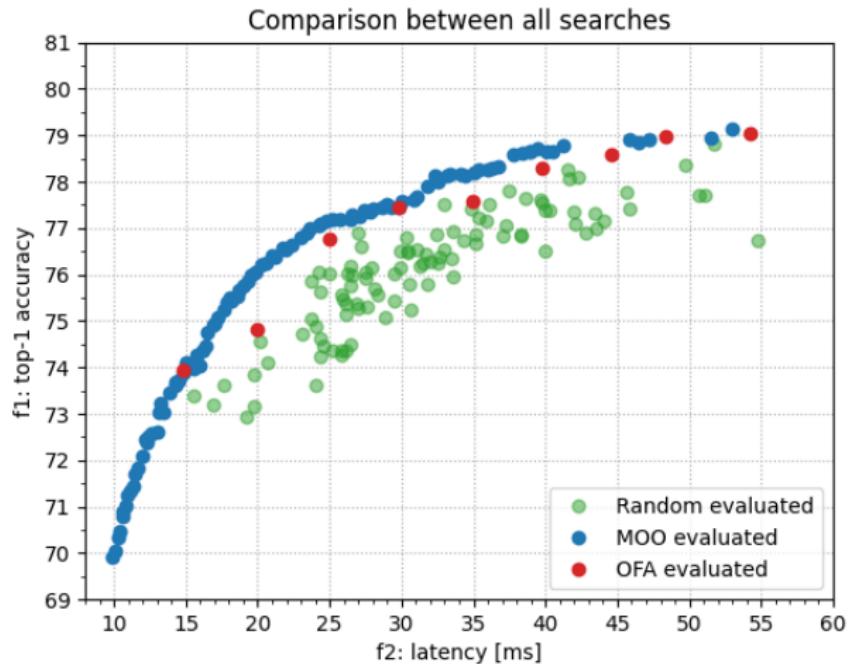
OFA<sup>2</sup> search (random)

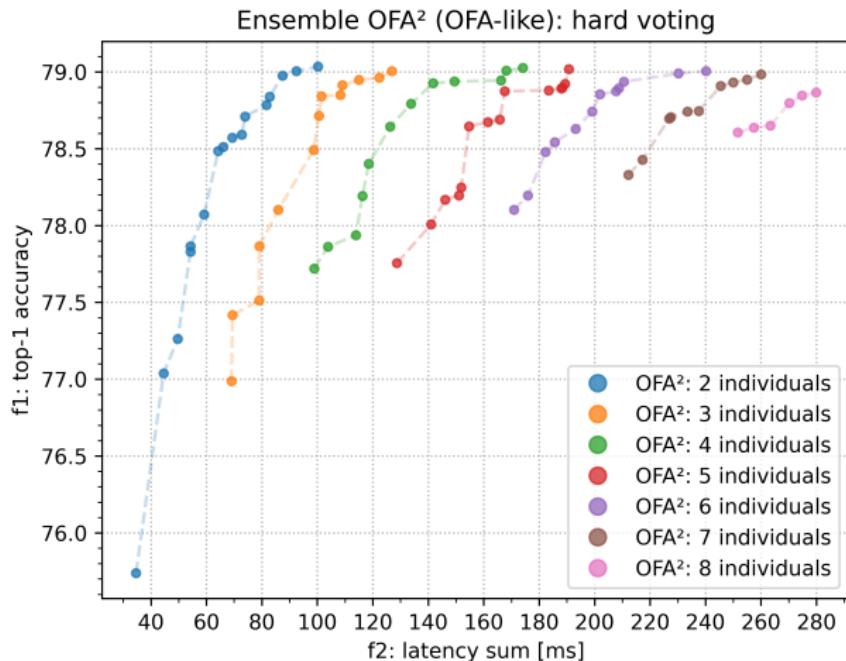
OFA<sup>2</sup> search (random): hard, sum

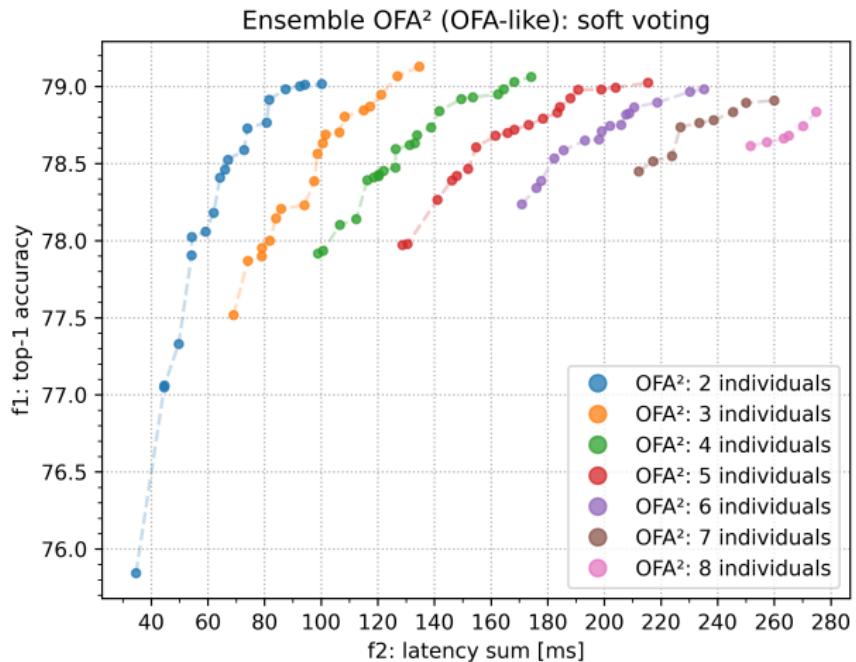
OFA<sup>2</sup> search (random): soft, sum

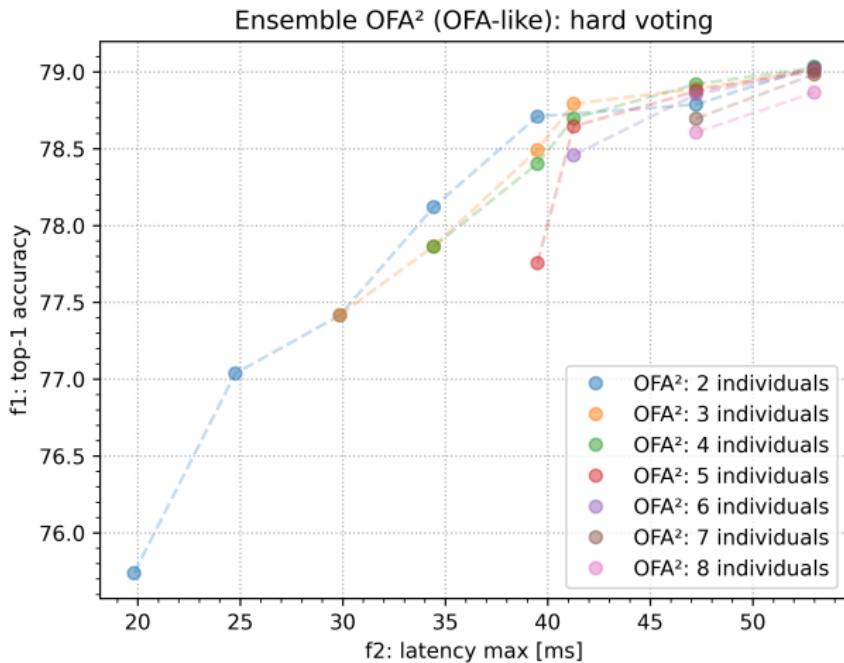
OFA<sup>2</sup> search (random): hard, max

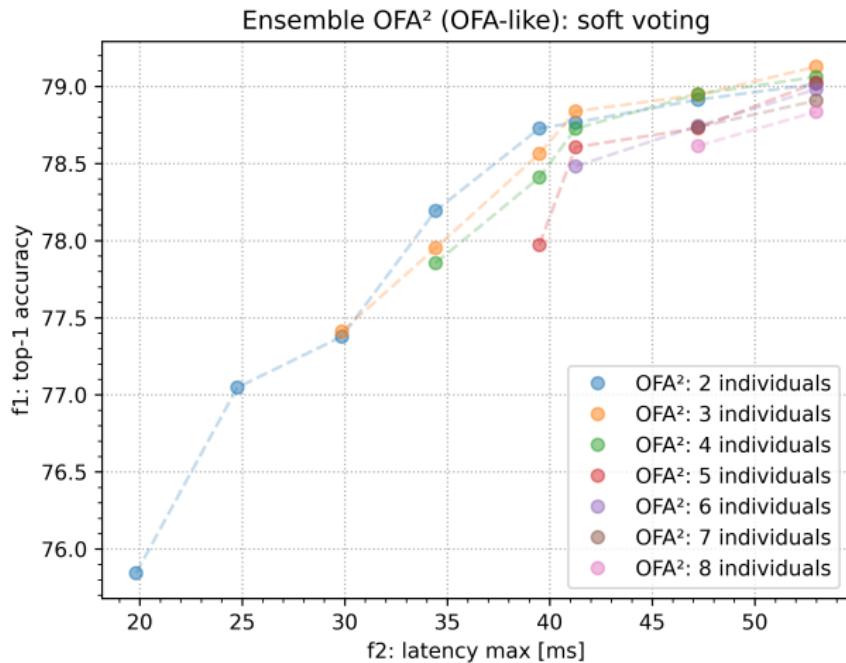
OFA<sup>2</sup> search (random): soft, sum

OFA<sup>2</sup> search (OFA-like)

OFA<sup>2</sup> search (OFA-like): hard, sum

OFA<sup>2</sup> search (OFA-like): soft, sum

OFA<sup>2</sup> search (OFA-like): hard, max

OFA<sup>2</sup> search (OFA-like): soft, sum

# OFA<sup>2</sup> x Random

- **Voting**

hard

soft

- **Latency**

sum

max

- **Individuals**

2 individuals

3 individuals

4 individuals

5 individuals

6 individuals

7 individuals

8 individuals

# OFA<sup>2</sup> x Random

- Voting

hard

soft

- Latency

sum

max

- Individuals

2 individuals

3 individuals

4 individuals

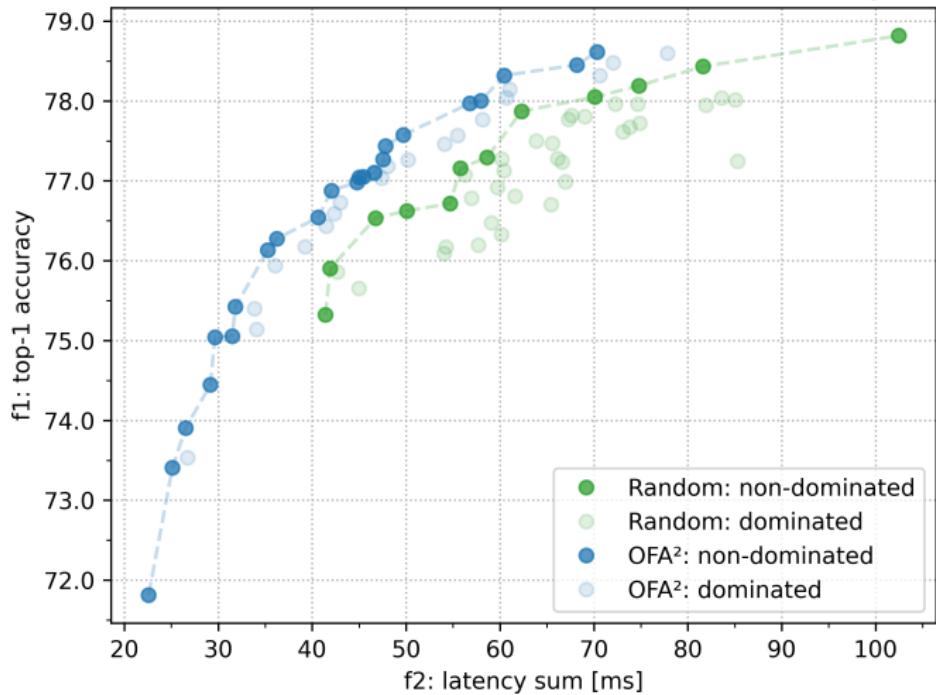
5 individuals

6 individuals

7 individuals

8 individuals

Ensemble Random x OFA<sup>2</sup>: 2 individuals (hard voting)



# OFA<sup>2</sup> x Random

- Voting

hard

soft

- Latency

sum

max

- Individuals

2 individuals

3 individuals

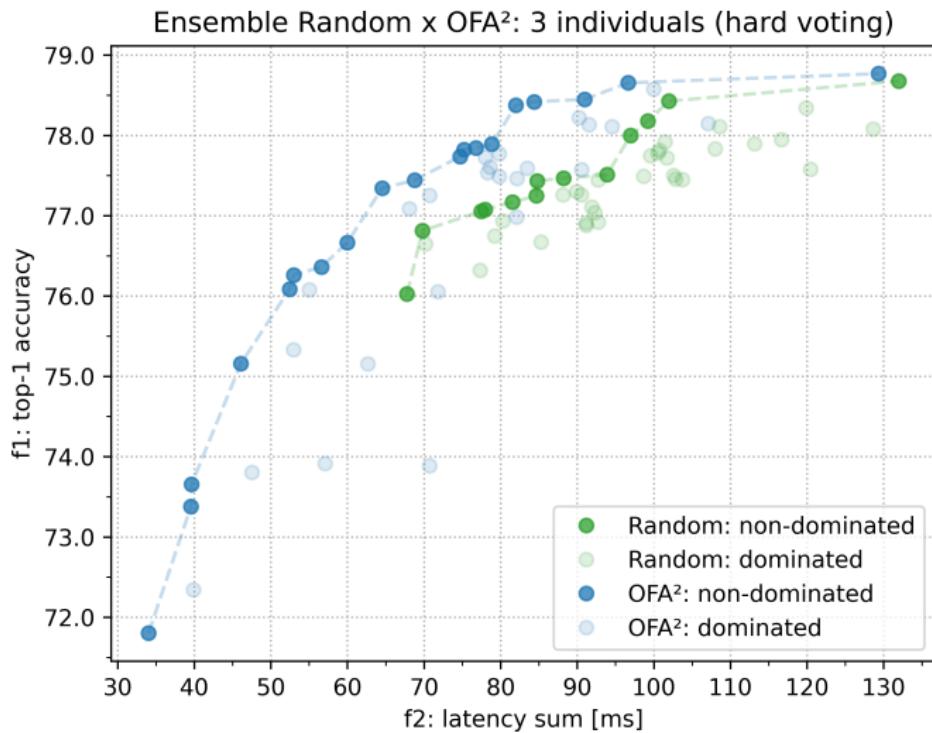
4 individuals

5 individuals

6 individuals

7 individuals

8 individuals



# OFA<sup>2</sup> x Random

- Voting

hard

soft

- Latency

sum

max

- Individuals

2 individuals

3 individuals

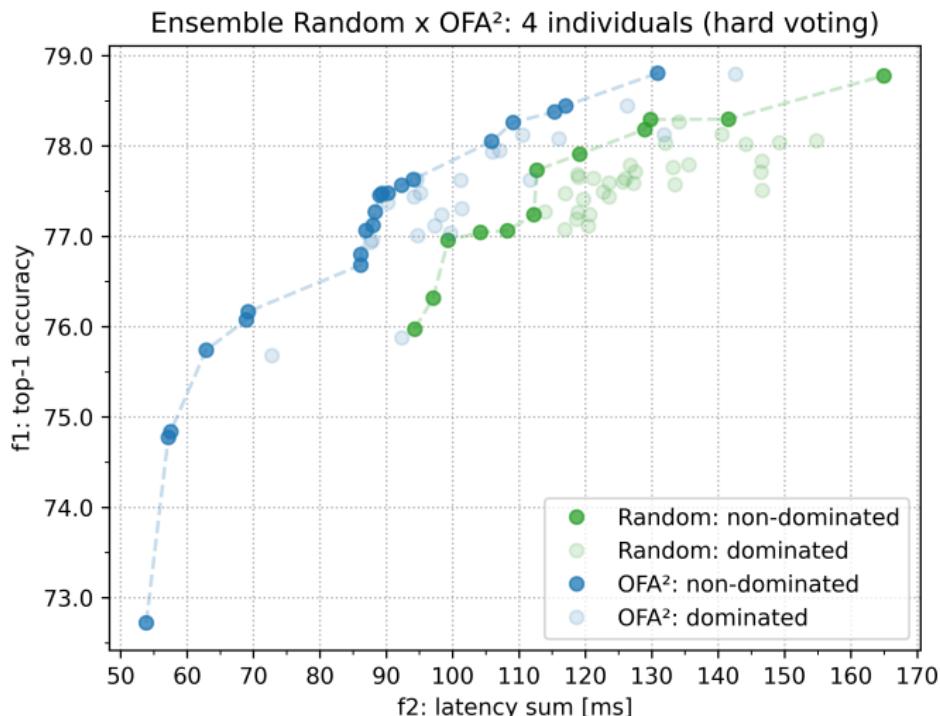
**4 individuals**

5 individuals

6 individuals

7 individuals

8 individuals



# OFA<sup>2</sup> x Random

- Voting

hard

soft

- Latency

sum

max

- Individuals

2 individuals

3 individuals

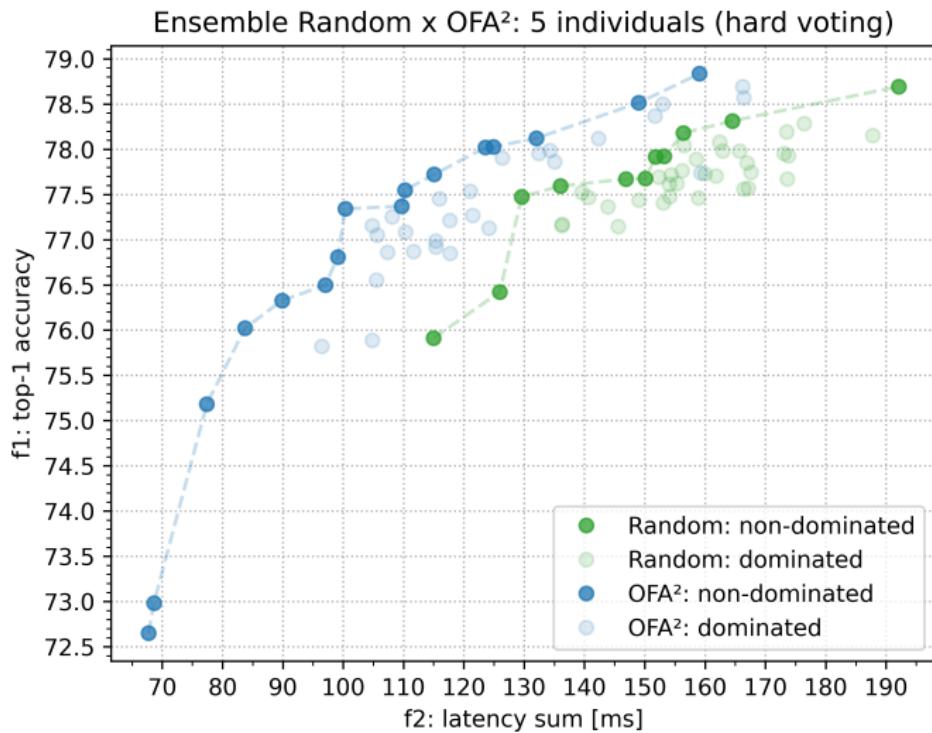
4 individuals

**5 individuals**

6 individuals

7 individuals

8 individuals



# OFA<sup>2</sup> x Random

- Voting

hard

soft

- Latency

sum

max

- Individuals

2 individuals

3 individuals

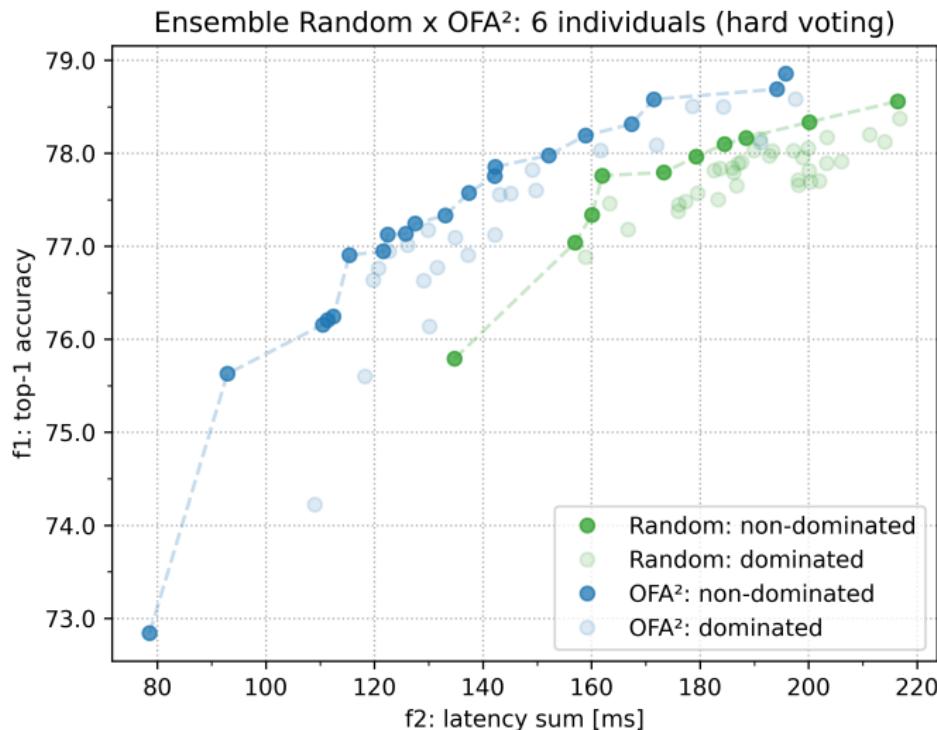
4 individuals

5 individuals

**6 individuals**

7 individuals

8 individuals



# OFA<sup>2</sup> x Random

- Voting

hard

soft

- Latency

sum

max

- Individuals

2 individuals

3 individuals

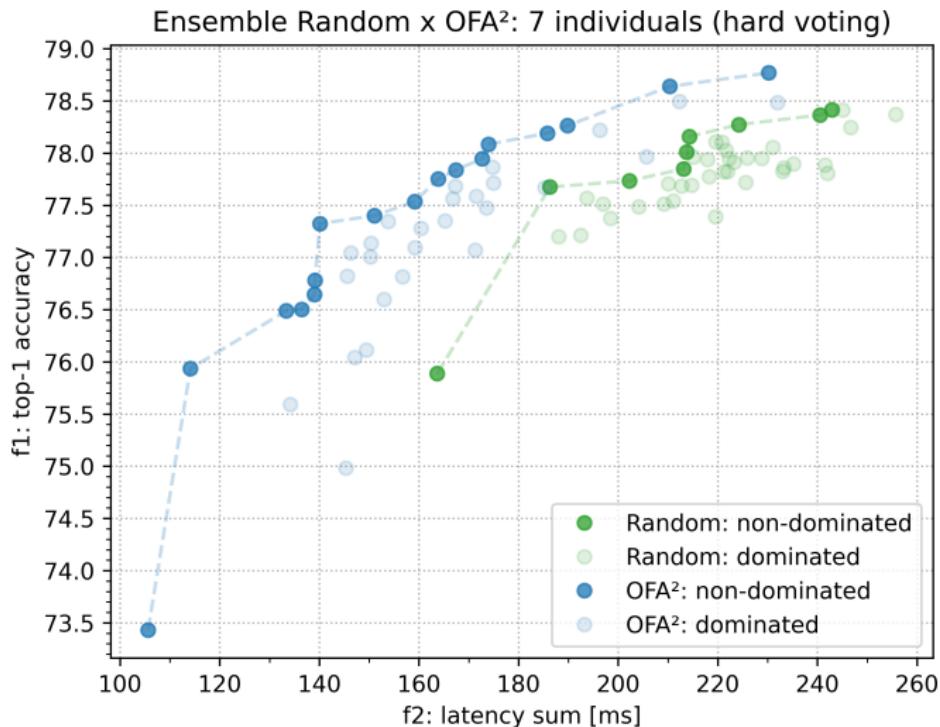
4 individuals

5 individuals

6 individuals

7 individuals

8 individuals



# OFA<sup>2</sup> x Random

- Voting

hard

soft

- Latency

sum

max

- Individuals

2 individuals

3 individuals

4 individuals

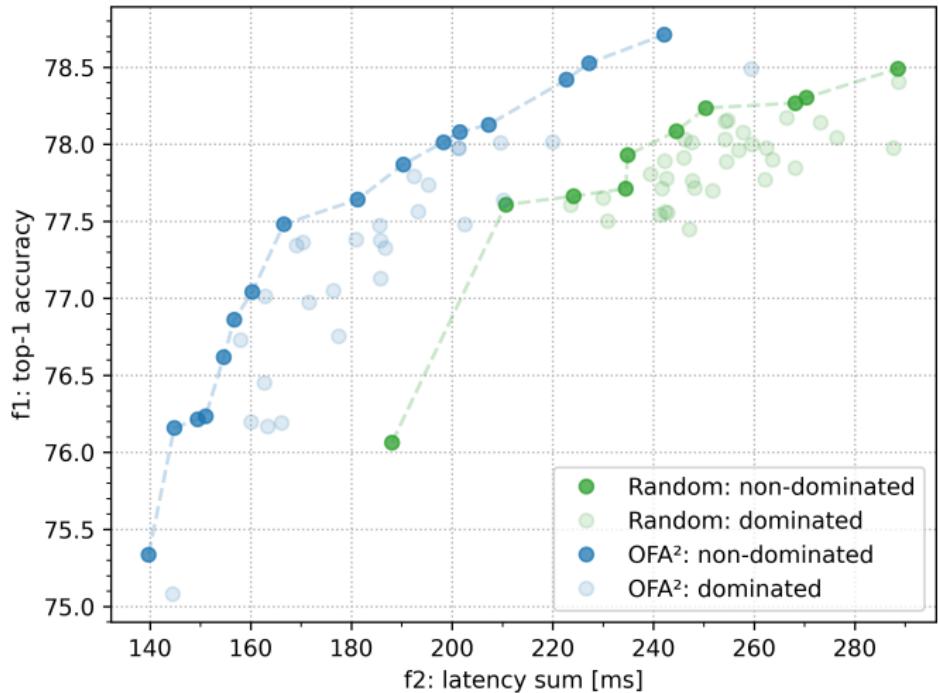
5 individuals

6 individuals

7 individuals

8 individuals

Ensemble Random x OFA<sup>2</sup>: 8 individuals (hard voting)



# OFA<sup>2</sup> x Random

- **Voting**

hard

soft

- **Latency**

sum

max

- **Individuals**

2 individuals

3 individuals

4 individuals

5 individuals

6 individuals

7 individuals

8 individuals

# OFA<sup>2</sup> x Random

- Voting

hard

soft

- Latency

sum

max

- Individuals

2 individuals

3 individuals

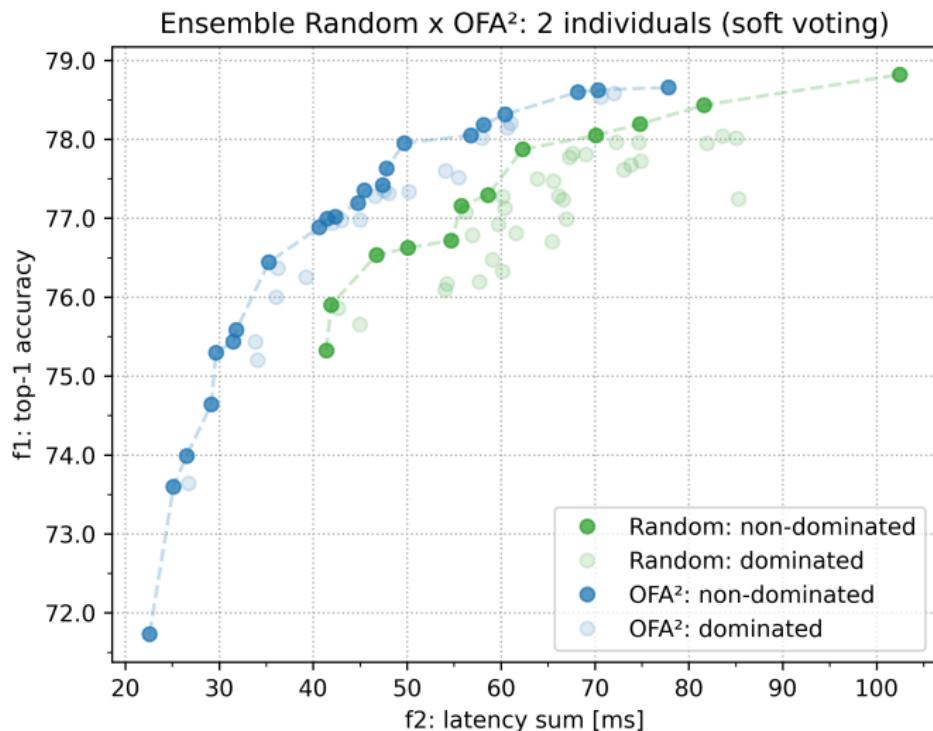
4 individuals

5 individuals

6 individuals

7 individuals

8 individuals



# OFA<sup>2</sup> x Random

- Voting

hard

soft

- Latency

sum

max

- Individuals

2 individuals

3 individuals

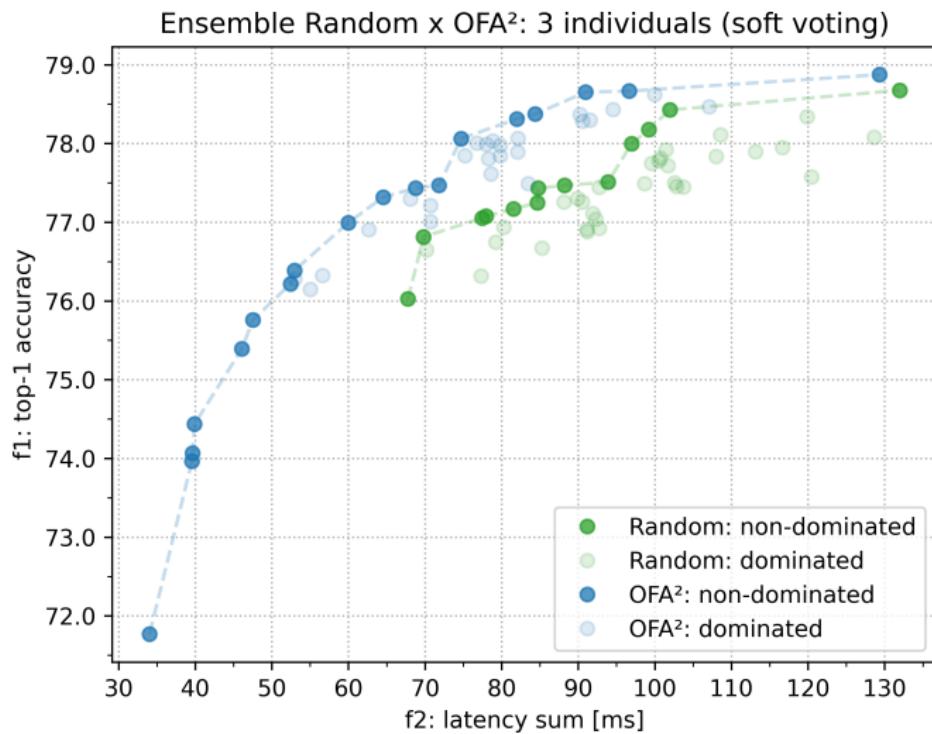
4 individuals

5 individuals

6 individuals

7 individuals

8 individuals



# OFA<sup>2</sup> x Random

- Voting

hard

soft

- Latency

sum

max

- Individuals

2 individuals

3 individuals

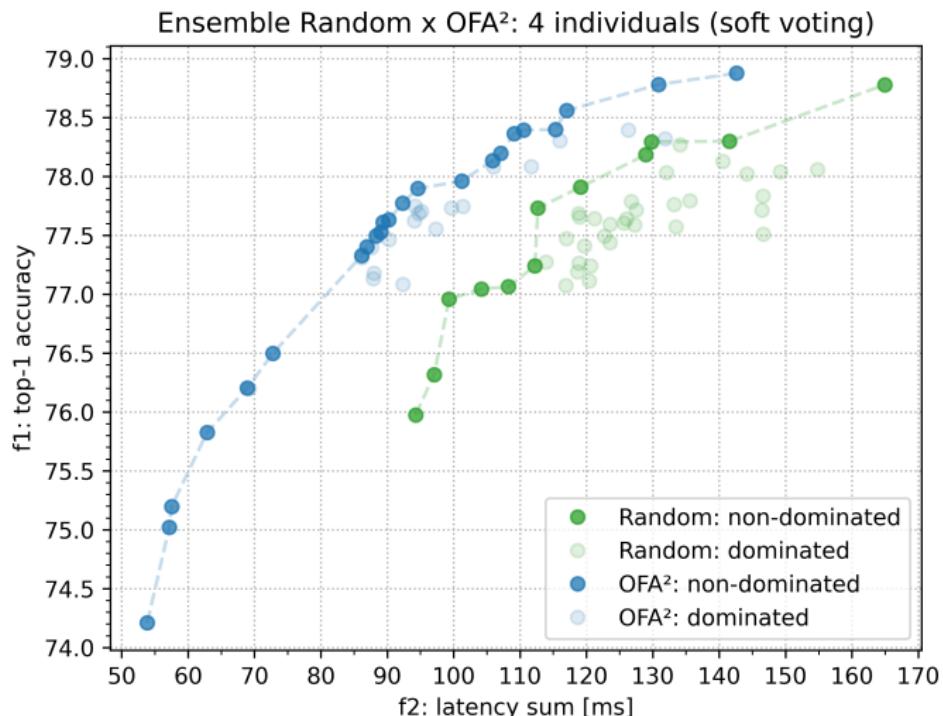
**4 individuals**

5 individuals

6 individuals

7 individuals

8 individuals



# OFA<sup>2</sup> x Random

- Voting

hard

soft

- Latency

sum

max

- Individuals

2 individuals

3 individuals

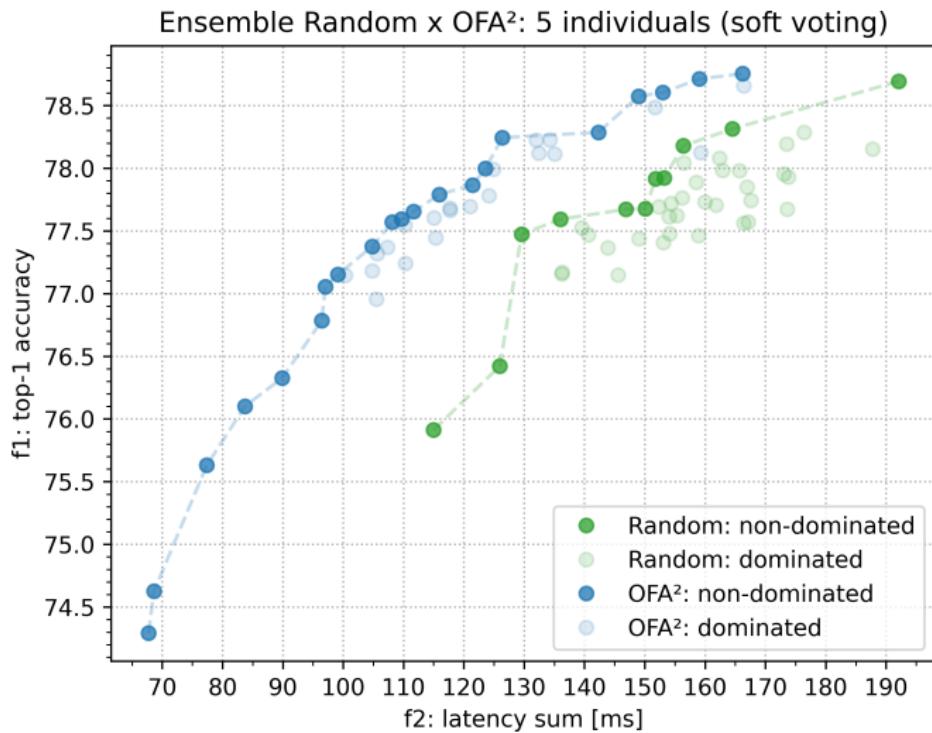
4 individuals

5 individuals

6 individuals

7 individuals

8 individuals



# OFA<sup>2</sup> x Random

- Voting

hard

soft

- Latency

sum

max

- Individuals

2 individuals

3 individuals

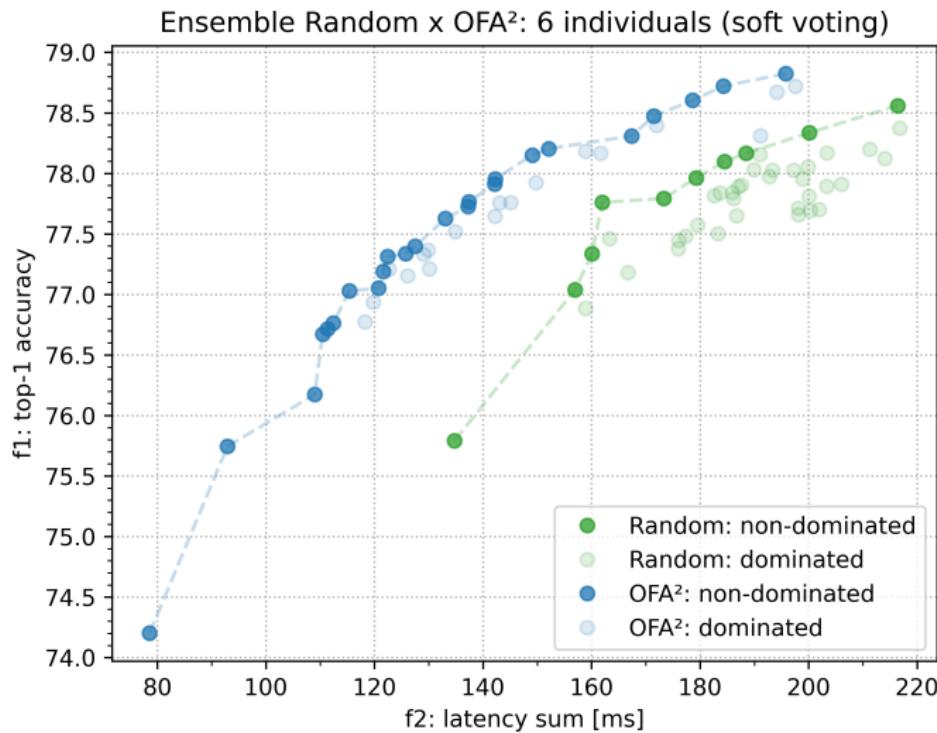
4 individuals

5 individuals

6 individuals

7 individuals

8 individuals



# OFA<sup>2</sup> x Random

- Voting

hard

soft

- Latency

sum

max

- Individuals

2 individuals

3 individuals

4 individuals

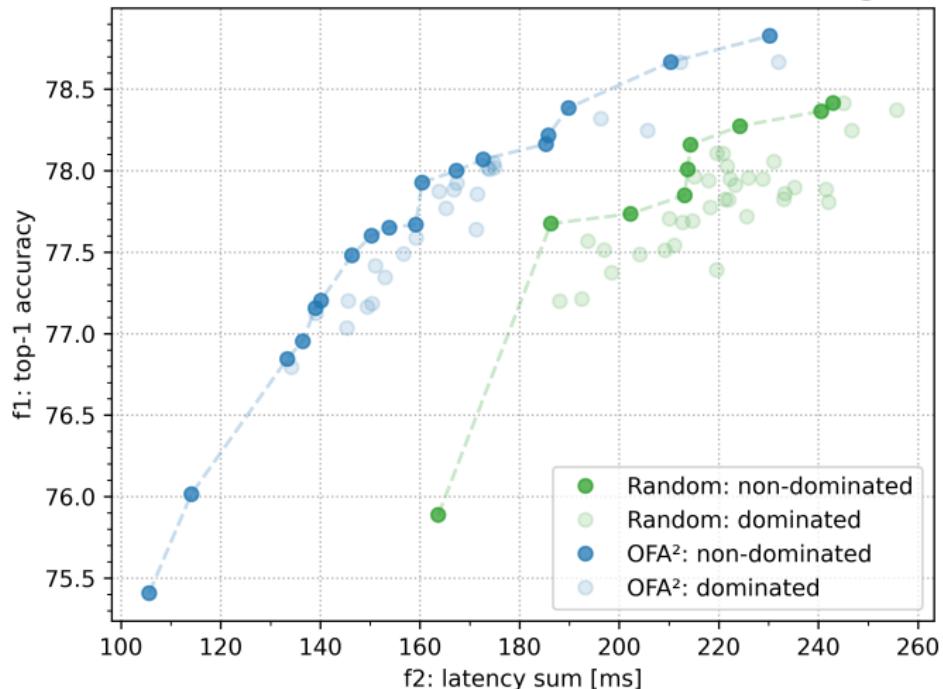
5 individuals

6 individuals

7 individuals

8 individuals

Ensemble Random x OFA<sup>2</sup>: 7 individuals (soft voting)



## OFA<sup>2</sup> x Random

- Voting

hard

soft

- Latency

sum

max

- Individuals

## 2 individuals

### 3 individuals

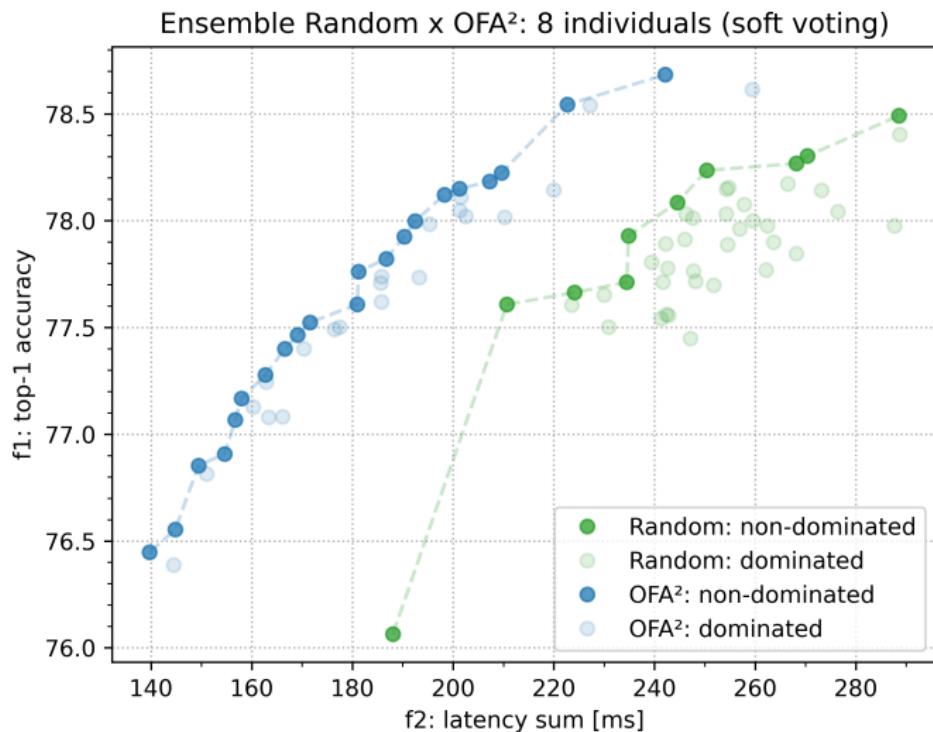
## 4 individuals

## 5 individuals

6 individuals

7 individuals

8 individuals



# OFA<sup>2</sup> x Random

- **Voting**

hard

soft

- **Latency**

sum

max

- **Individuals**

2 individuals

3 individuals

4 individuals

5 individuals

6 individuals

7 individuals

8 individuals

# OFA<sup>2</sup> x Random

- Voting

hard

soft

- Latency

sum

max

- Individuals

2 individuals

3 individuals

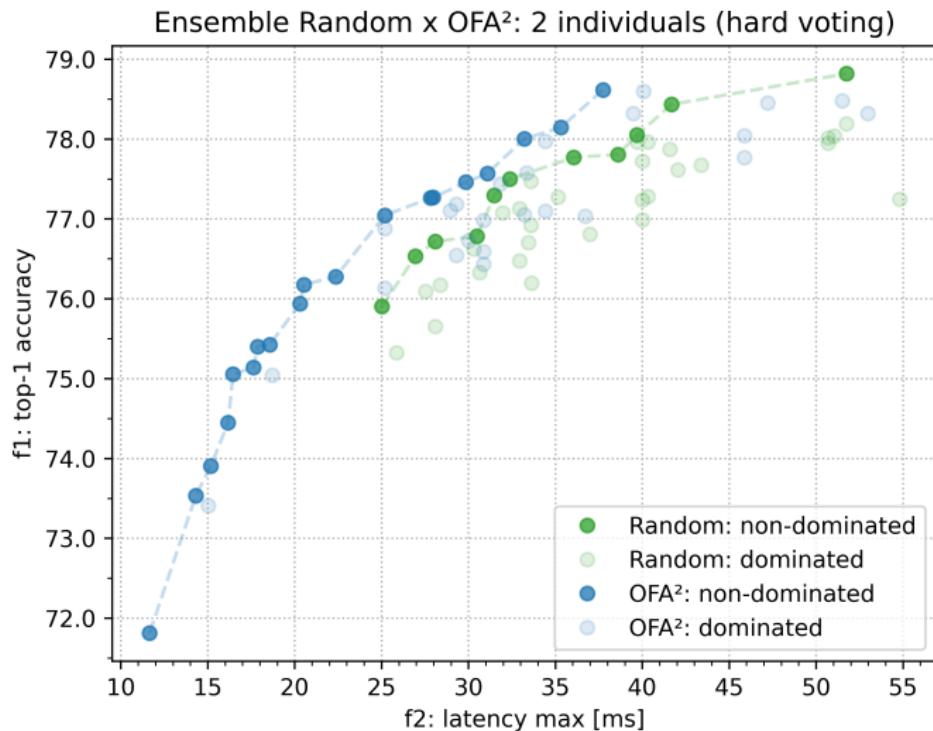
4 individuals

5 individuals

6 individuals

7 individuals

8 individuals



# OFA<sup>2</sup> x Random

- Voting

hard

soft

- Latency

sum

max

- Individuals

2 individuals

3 individuals

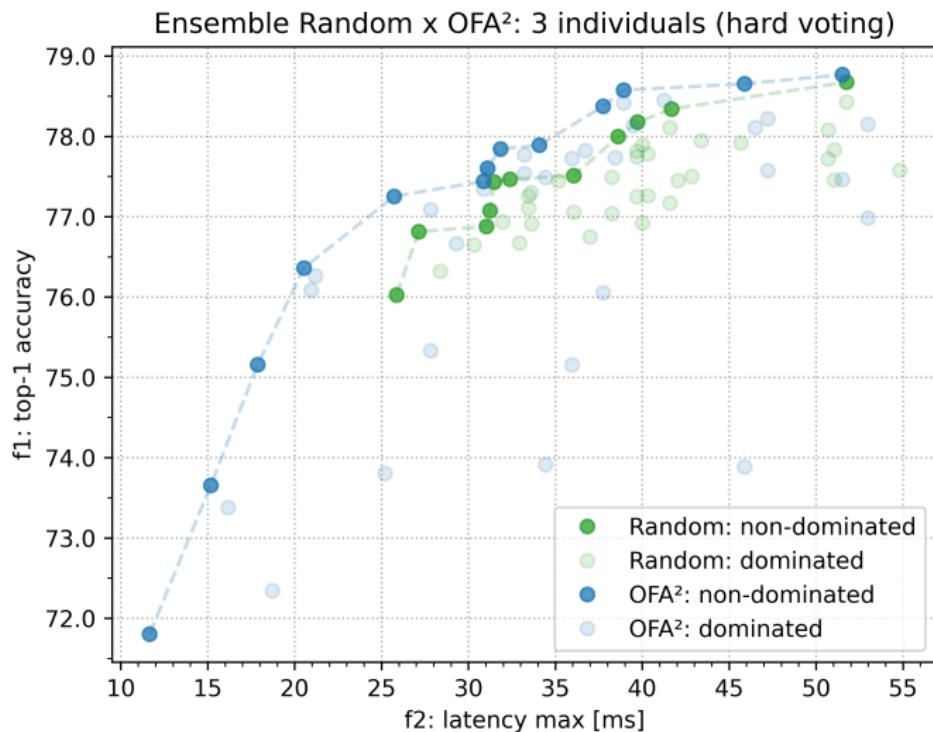
4 individuals

5 individuals

6 individuals

7 individuals

8 individuals



# OFA<sup>2</sup> x Random

- Voting

hard

soft

- Latency

sum

max

- Individuals

2 individuals

3 individuals

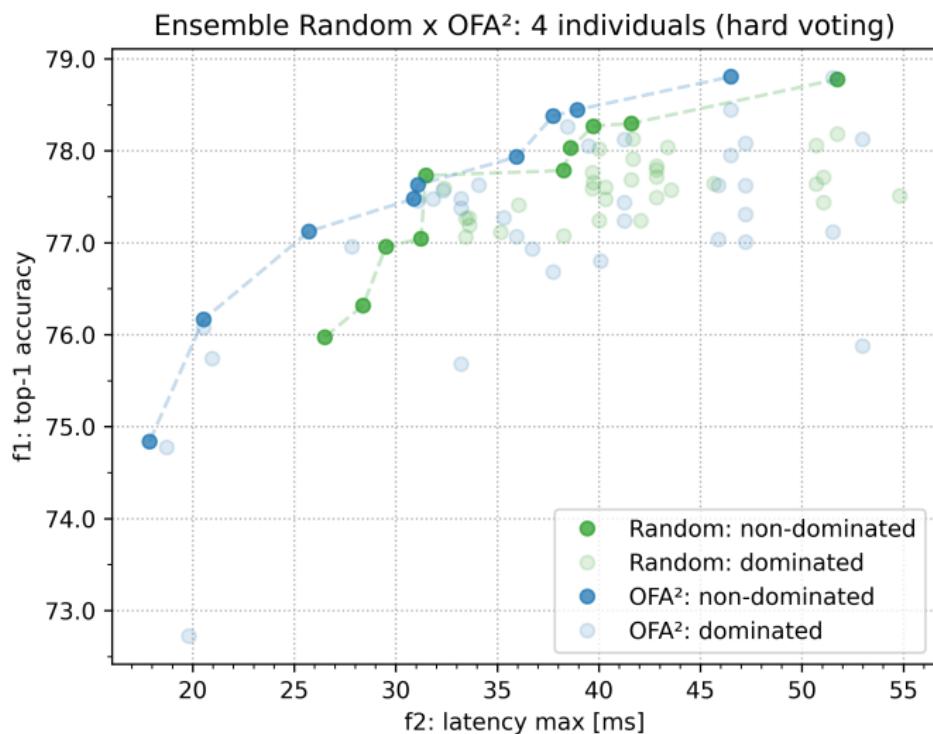
**4 individuals**

5 individuals

6 individuals

7 individuals

8 individuals



# OFA<sup>2</sup> x Random

- Voting

hard

soft

- Latency

sum

max

- Individuals

2 individuals

3 individuals

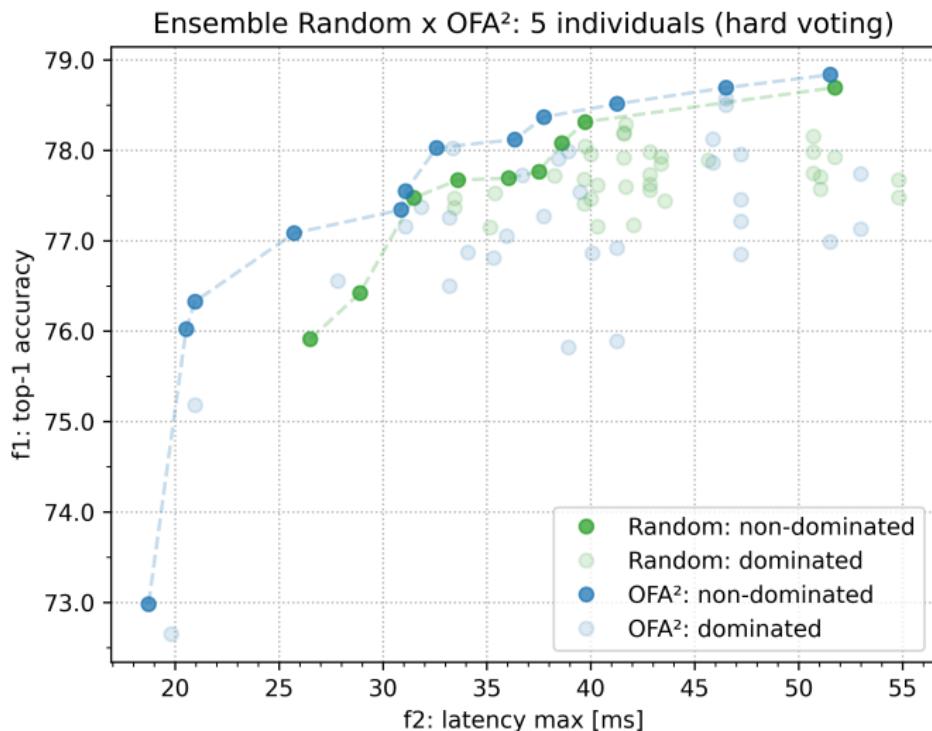
4 individuals

**5 individuals**

6 individuals

7 individuals

8 individuals



# OFA<sup>2</sup> x Random

- Voting

hard

soft

- Latency

sum

max

- Individuals

2 individuals

3 individuals

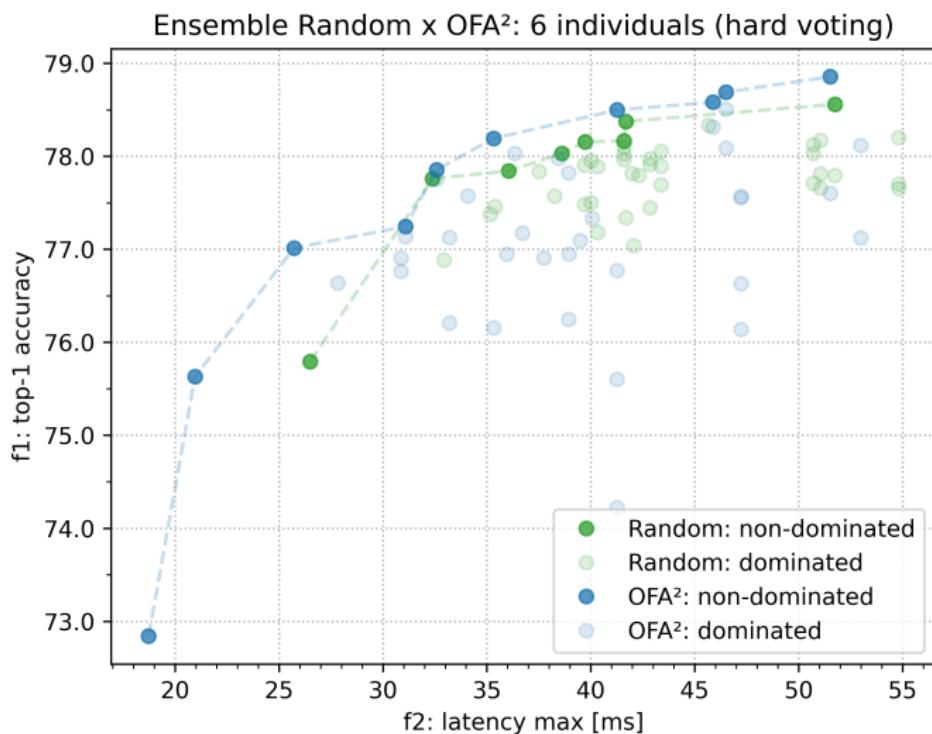
4 individuals

5 individuals

**6 individuals**

7 individuals

8 individuals



# OFA<sup>2</sup> x Random

- Voting

hard

soft

- Latency

sum

max

- Individuals

2 individuals

3 individuals

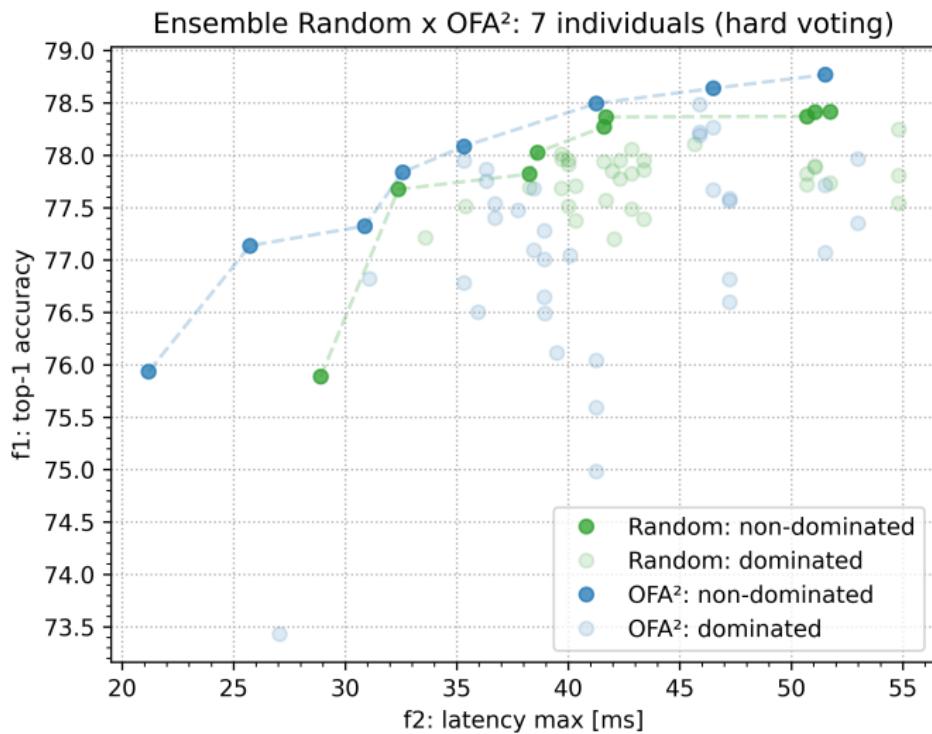
4 individuals

5 individuals

6 individuals

7 individuals

8 individuals



# OFA<sup>2</sup> x Random

- Voting

hard

soft

- Latency

sum

max

- Individuals

2 individuals

3 individuals

4 individuals

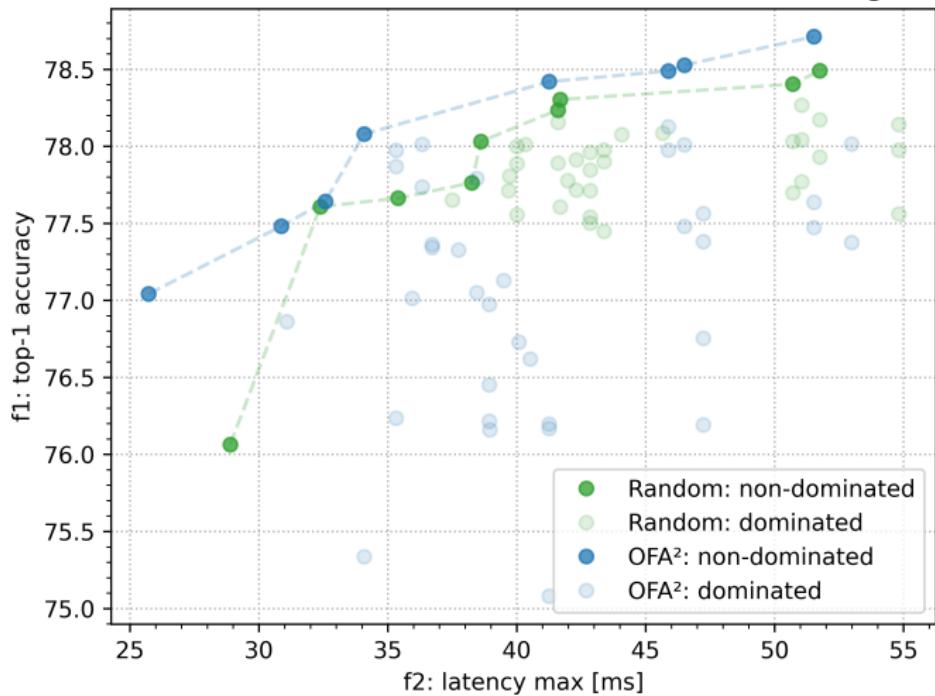
5 individuals

6 individuals

7 individuals

8 individuals

Ensemble Random x OFA<sup>2</sup>: 8 individuals (hard voting)



# OFA<sup>2</sup> x Random

- **Voting**

hard

soft

- **Latency**

sum

max

- **Individuals**

2 individuals

3 individuals

4 individuals

5 individuals

6 individuals

7 individuals

8 individuals

# OFA<sup>2</sup> x Random

- Voting

hard

soft

- Latency

sum

max

- Individuals

2 individuals

3 individuals

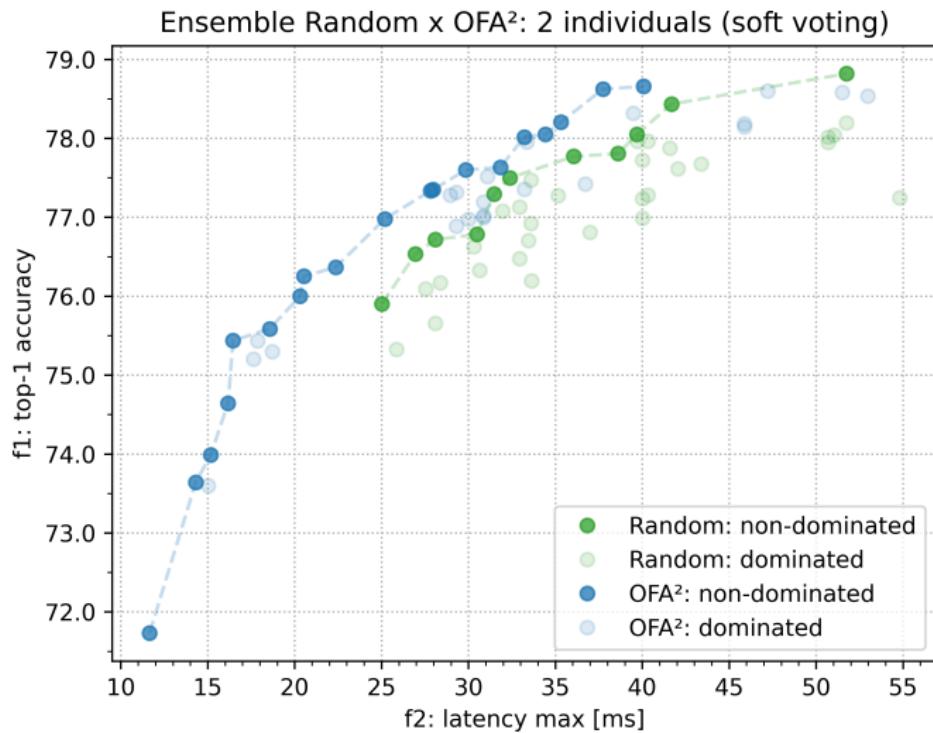
4 individuals

5 individuals

6 individuals

7 individuals

8 individuals



# OFA<sup>2</sup> x Random

- Voting

hard

**soft**

- Latency

sum

**max**

- Individuals

2 individuals

**3 individuals**

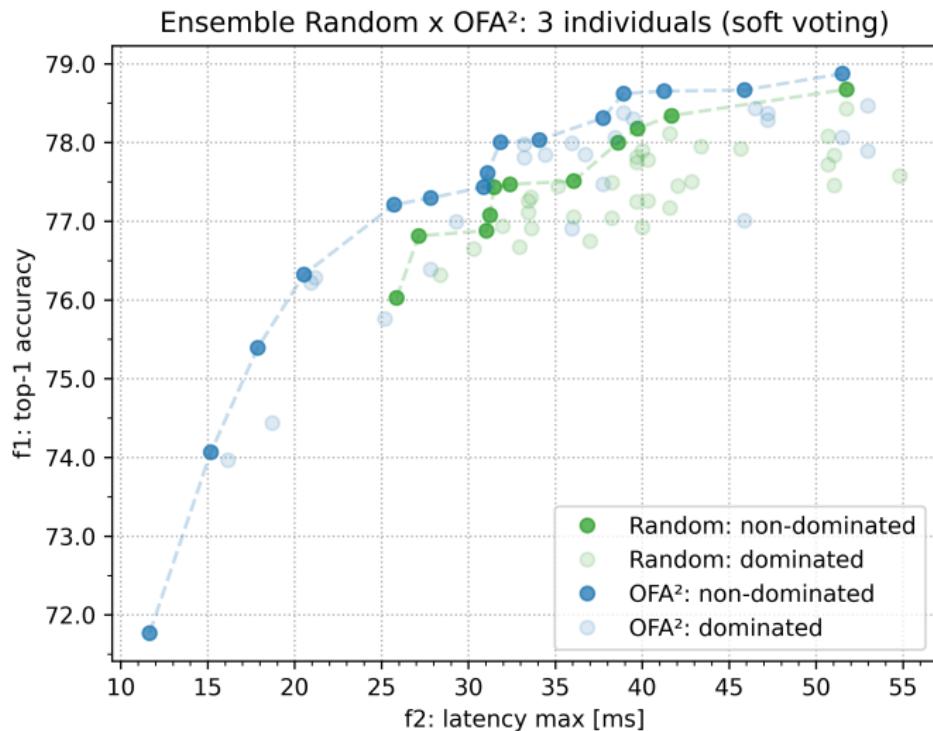
4 individuals

5 individuals

6 individuals

7 individuals

8 individuals



# OFA<sup>2</sup> x Random

- Voting

hard

soft

- Latency

sum

max

- Individuals

2 individuals

3 individuals

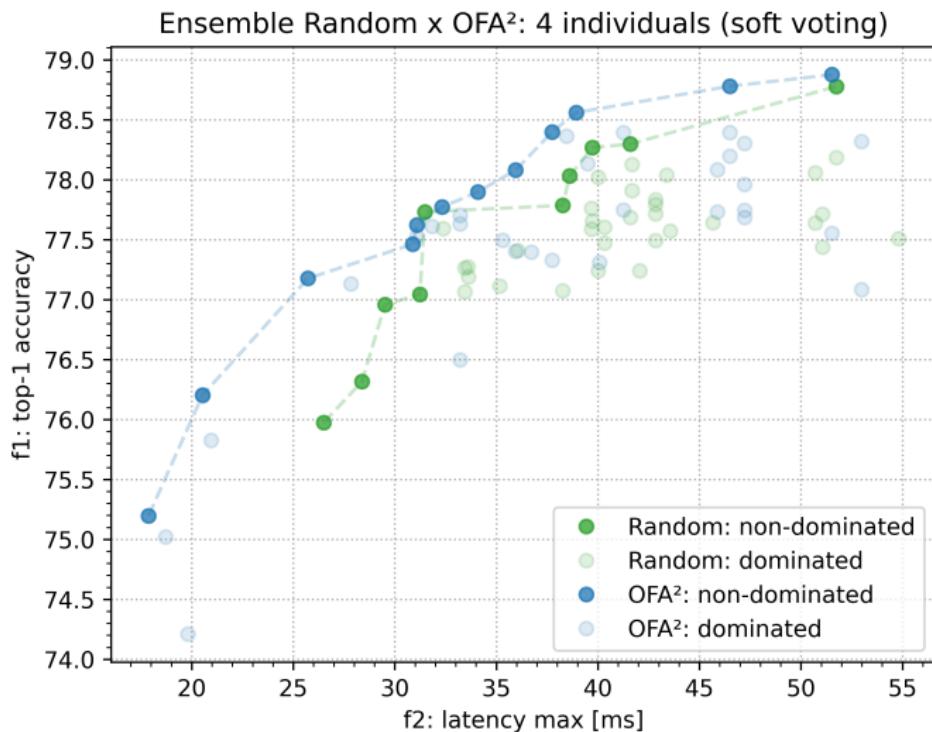
**4 individuals**

5 individuals

6 individuals

7 individuals

8 individuals



# OFA<sup>2</sup> x Random

- Voting

hard

**soft**

- Latency

sum

**max**

- Individuals

2 individuals

3 individuals

4 individuals

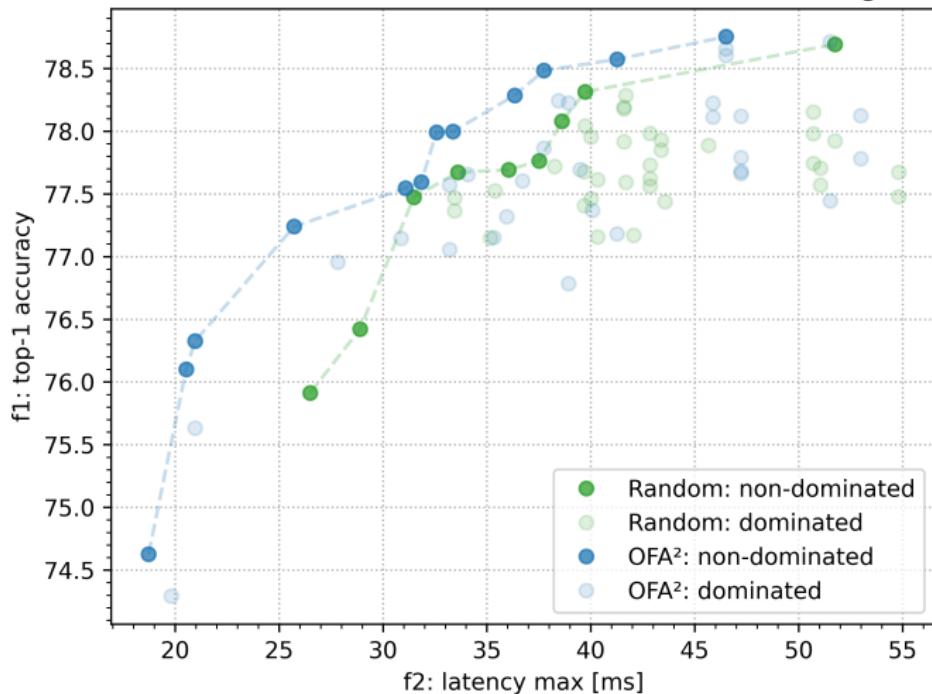
**5 individuals**

6 individuals

7 individuals

8 individuals

Ensemble Random x OFA<sup>2</sup>: 5 individuals (soft voting)



# OFA<sup>2</sup> x Random

- Voting

hard

soft

- Latency

sum

max

- Individuals

2 individuals

3 individuals

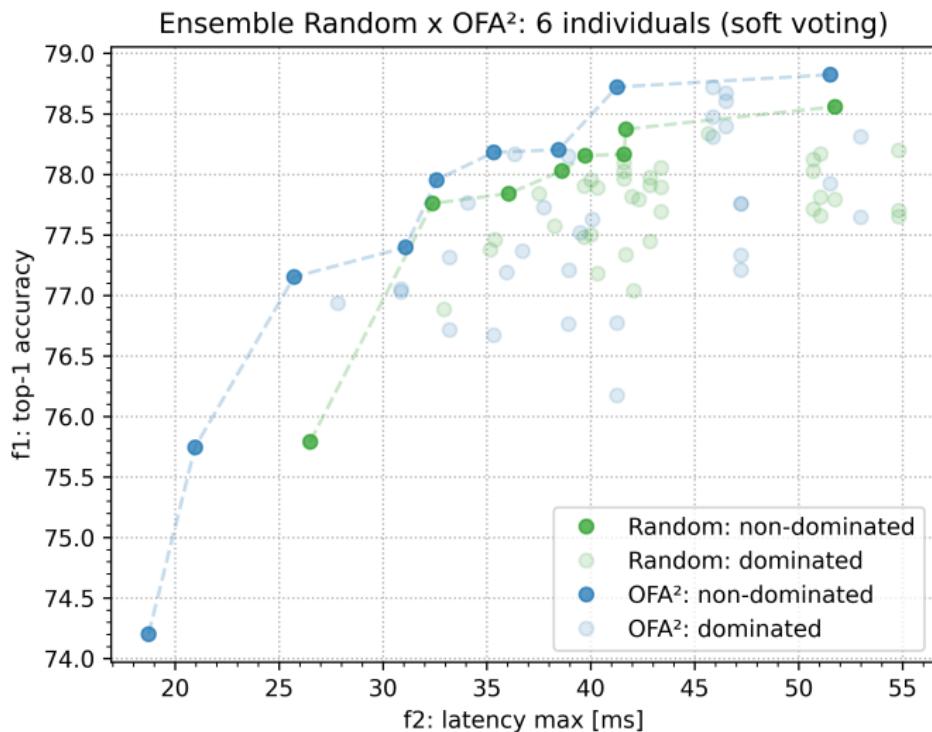
4 individuals

5 individuals

6 individuals

7 individuals

8 individuals



# OFA<sup>2</sup> x Random

- Voting

hard

soft

- Latency

sum

max

- Individuals

2 individuals

3 individuals

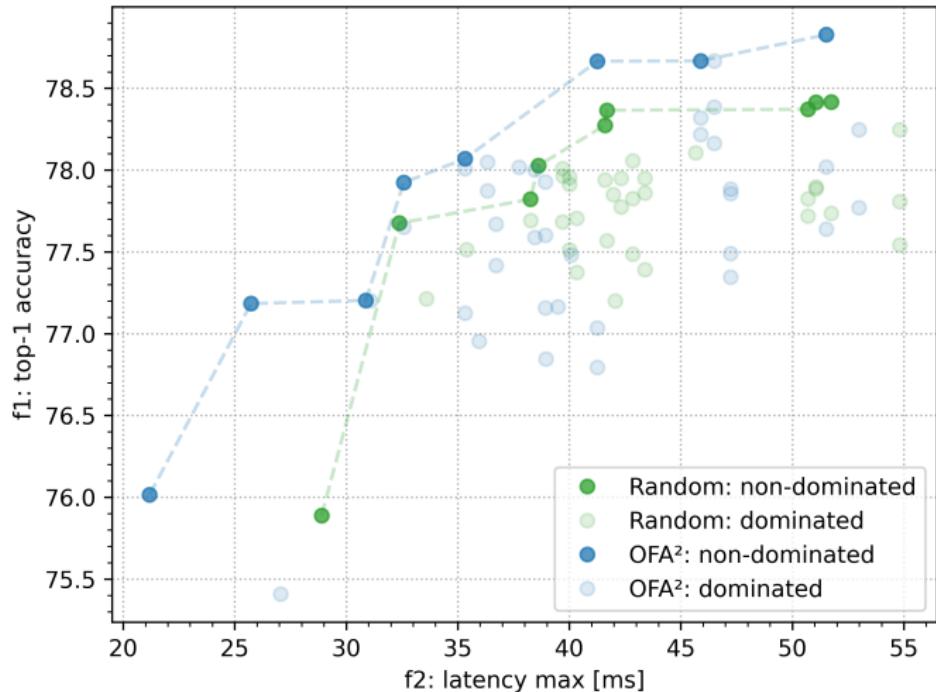
4 individuals

5 individuals

6 individuals

7 individuals

8 individuals

Ensemble Random x OFA<sup>2</sup>: 7 individuals (soft voting)

# OFA<sup>2</sup> x Random

- Voting

hard

soft

- Latency

sum

max

- Individuals

2 individuals

3 individuals

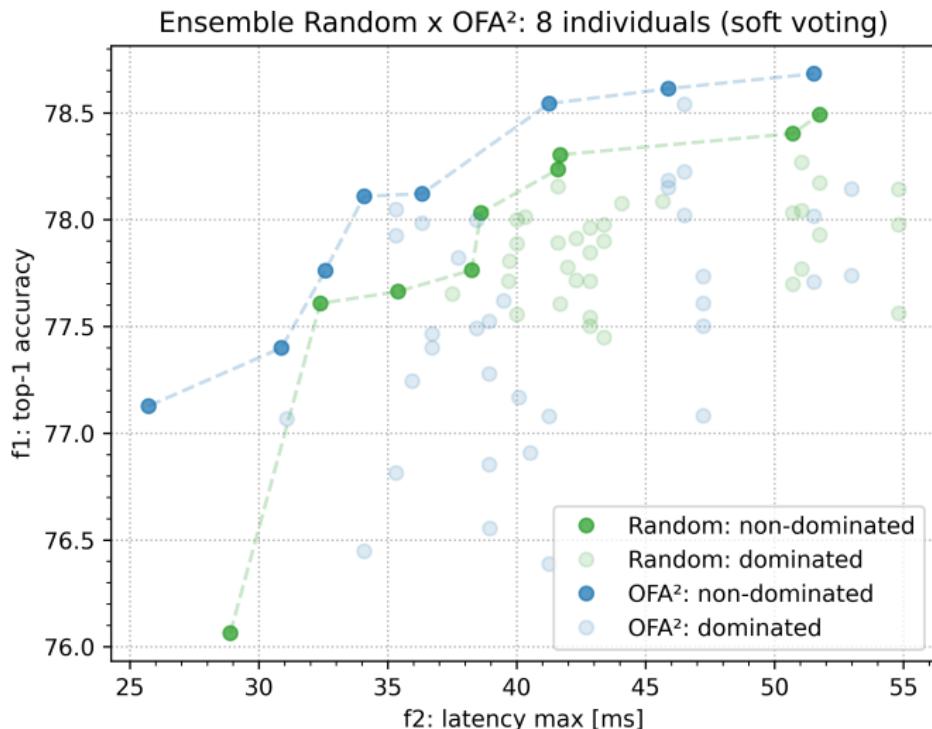
4 individuals

5 individuals

6 individuals

7 individuals

8 individuals



# OFA<sup>2</sup> × OFA

- **Voting**

hard

soft

- **Latency**

sum

max

- **Individuals**

2 individuals

3 individuals

4 individuals

5 individuals

6 individuals

7 individuals

8 individuals

# OFA<sup>2</sup> x OFA

- Voting

hard

soft

- Latency

sum

max

- Individuals

2 individuals

3 individuals

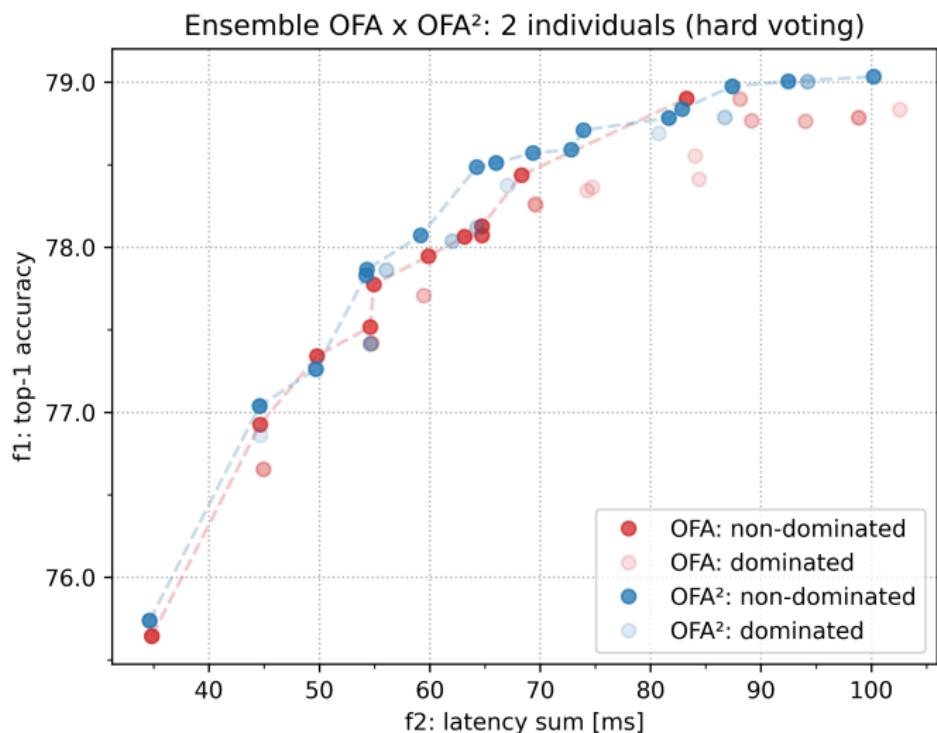
4 individuals

5 individuals

6 individuals

7 individuals

8 individuals



# OFA<sup>2</sup> x OFA

- Voting

hard

soft

- Latency

sum

max

- Individuals

2 individuals

**3 individuals**

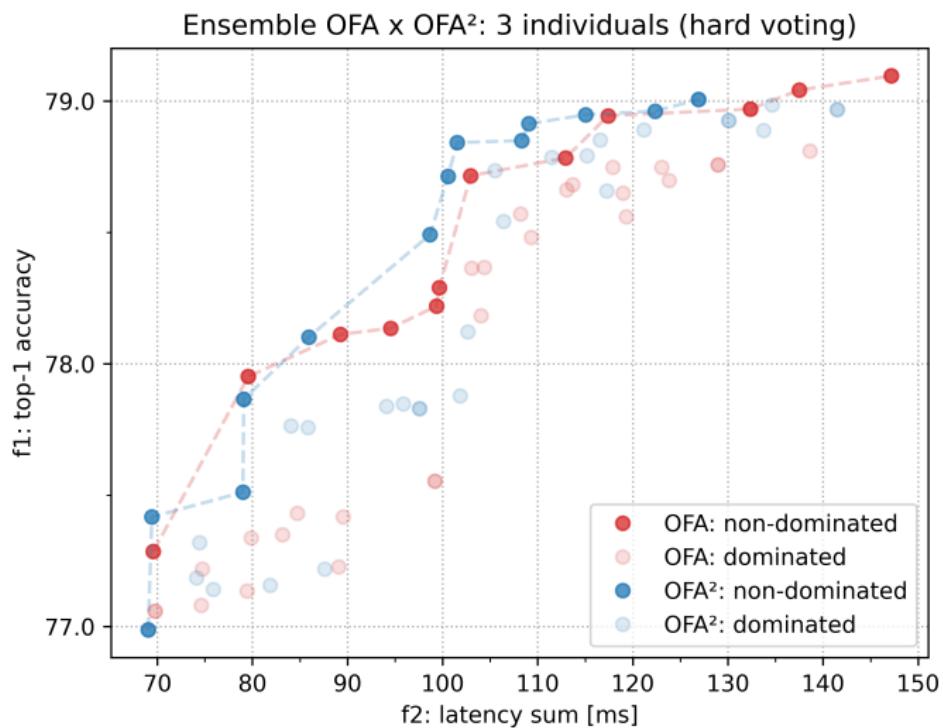
4 individuals

5 individuals

6 individuals

7 individuals

8 individuals



# OFA<sup>2</sup> x OFA

- Voting

hard

soft

- Latency

sum

max

- Individuals

2 individuals

3 individuals

**4 individuals**

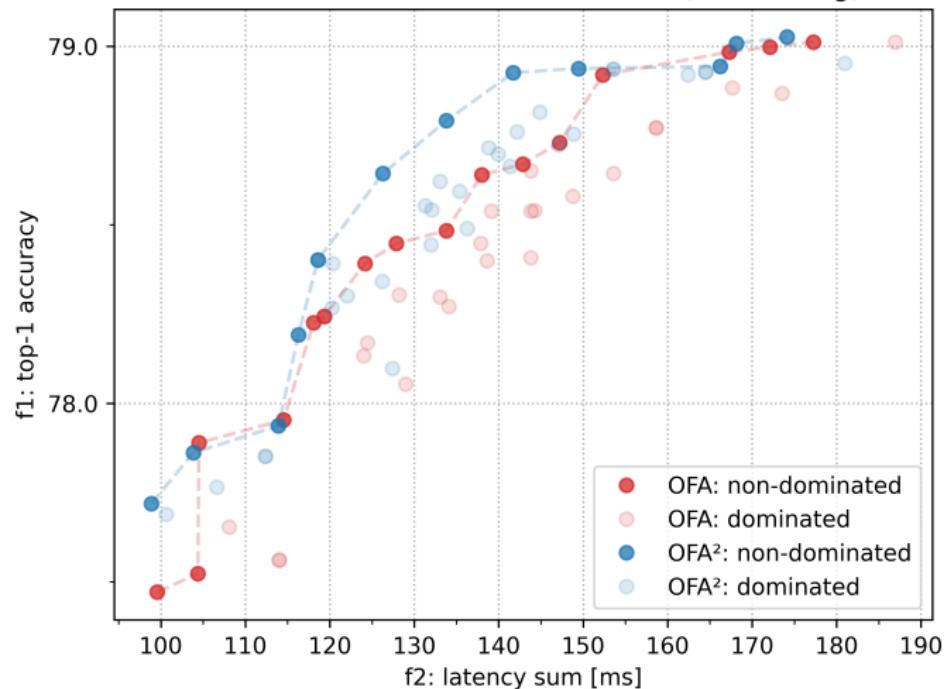
5 individuals

6 individuals

7 individuals

8 individuals

Ensemble OFA x OFA<sup>2</sup>: 4 individuals (hard voting)



# OFA<sup>2</sup> x OFA

- Voting

hard

soft

- Latency

sum

max

- Individuals

2 individuals

3 individuals

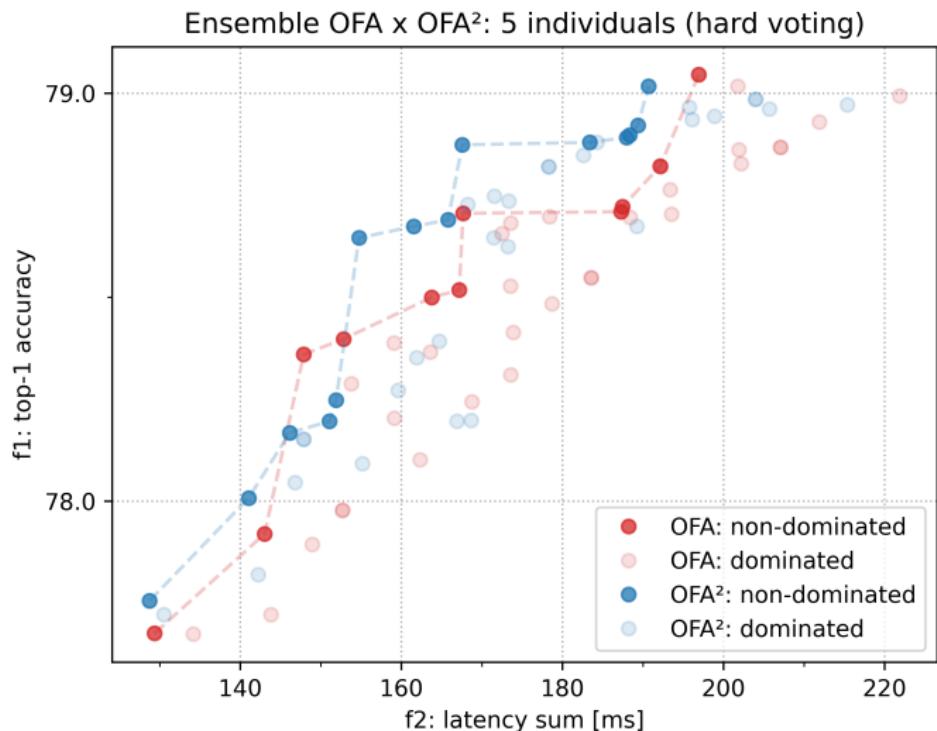
4 individuals

**5 individuals**

6 individuals

7 individuals

8 individuals



# OFA<sup>2</sup> x OFA

- Voting

hard

soft

- Latency

sum

max

- Individuals

2 individuals

3 individuals

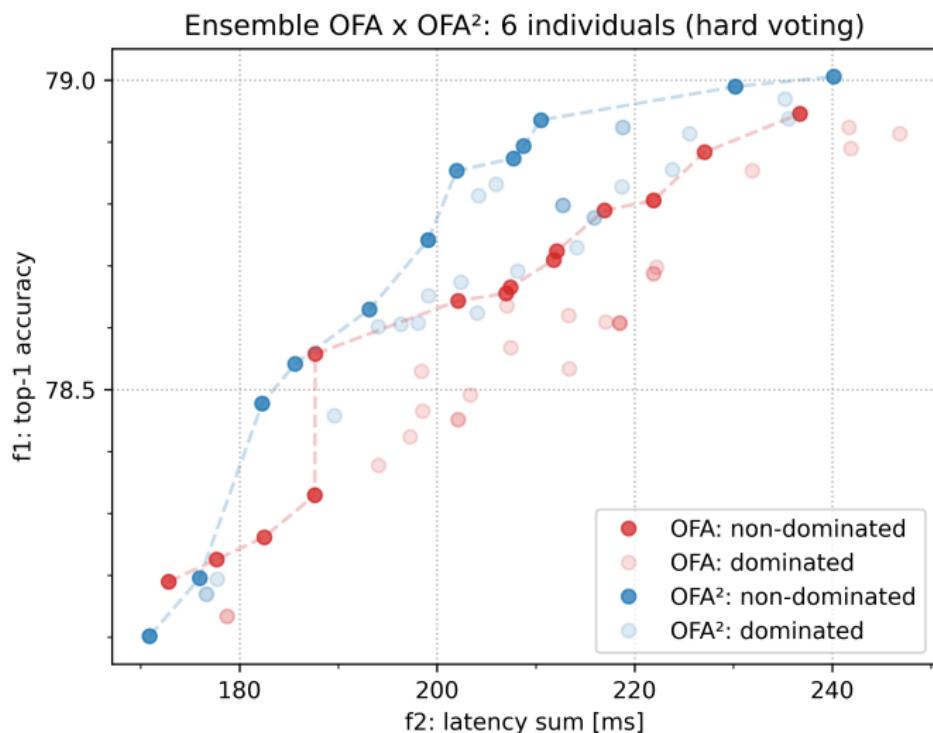
4 individuals

5 individuals

6 individuals

7 individuals

8 individuals



# OFA<sup>2</sup> x OFA

- Voting

hard

soft

- Latency

sum

max

- Individuals

2 individuals

3 individuals

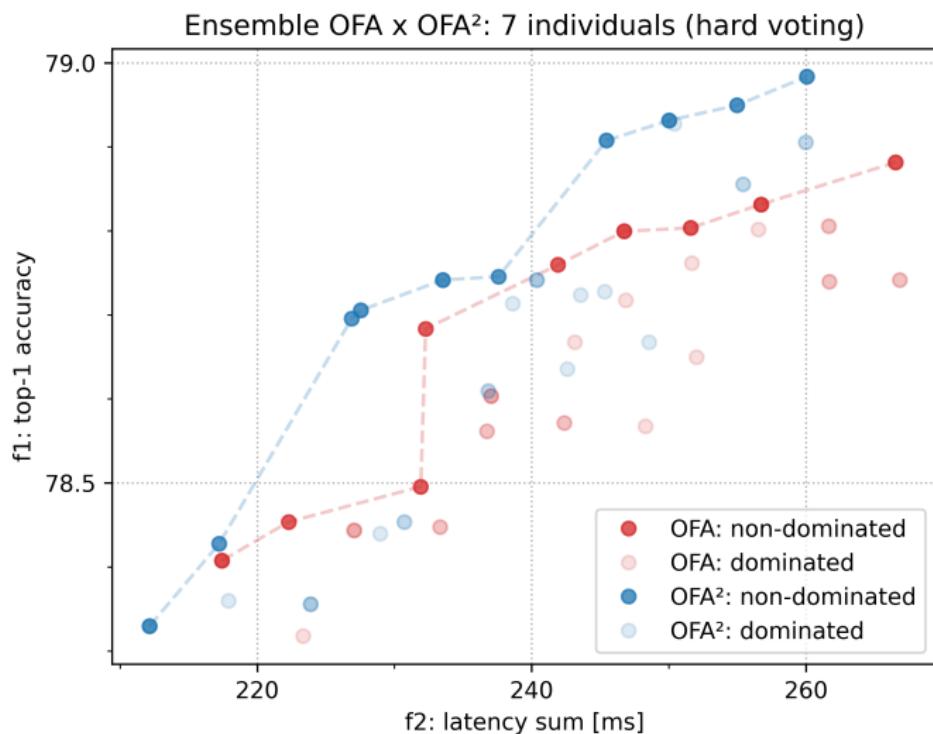
4 individuals

5 individuals

6 individuals

7 individuals

8 individuals



# OFA<sup>2</sup> x OFA

- Voting

hard

soft

- Latency

sum

max

- Individuals

2 individuals

3 individuals

4 individuals

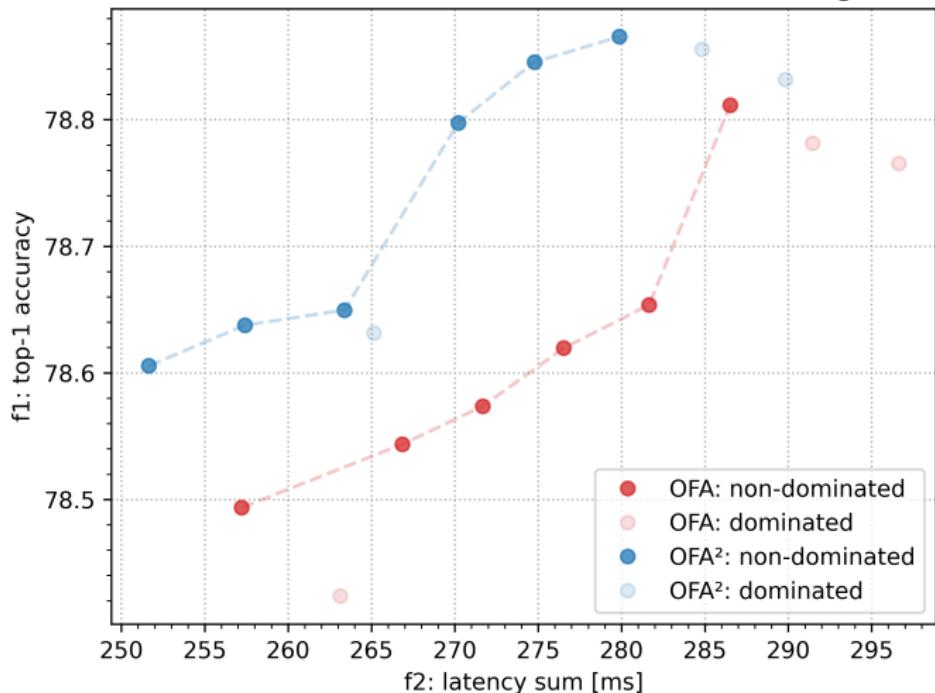
5 individuals

6 individuals

7 individuals

8 individuals

Ensemble OFA x OFA<sup>2</sup>: 8 individuals (hard voting)



# OFA<sup>2</sup> × OFA

- **Voting**

hard

soft

- **Latency**

sum

max

- **Individuals**

2 individuals

3 individuals

4 individuals

5 individuals

6 individuals

7 individuals

8 individuals

# OFA<sup>2</sup> x OFA

- Voting

hard

soft

- Latency

sum

max

- Individuals

2 individuals

3 individuals

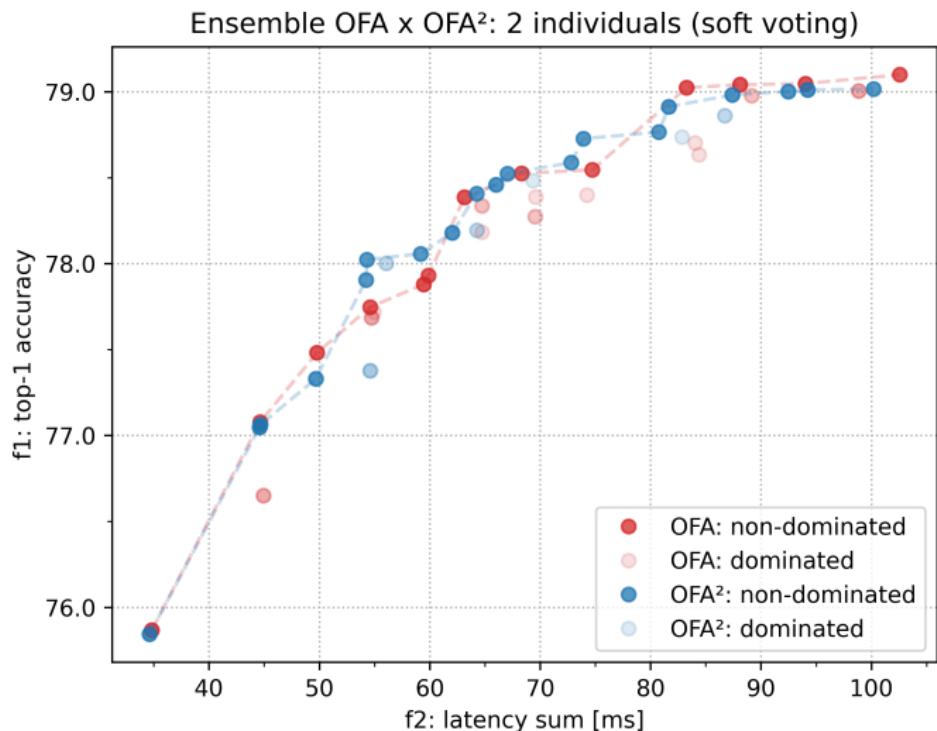
4 individuals

5 individuals

6 individuals

7 individuals

8 individuals



# OFA<sup>2</sup> x OFA

- Voting

hard

soft

- Latency

sum

max

- Individuals

2 individuals

3 individuals

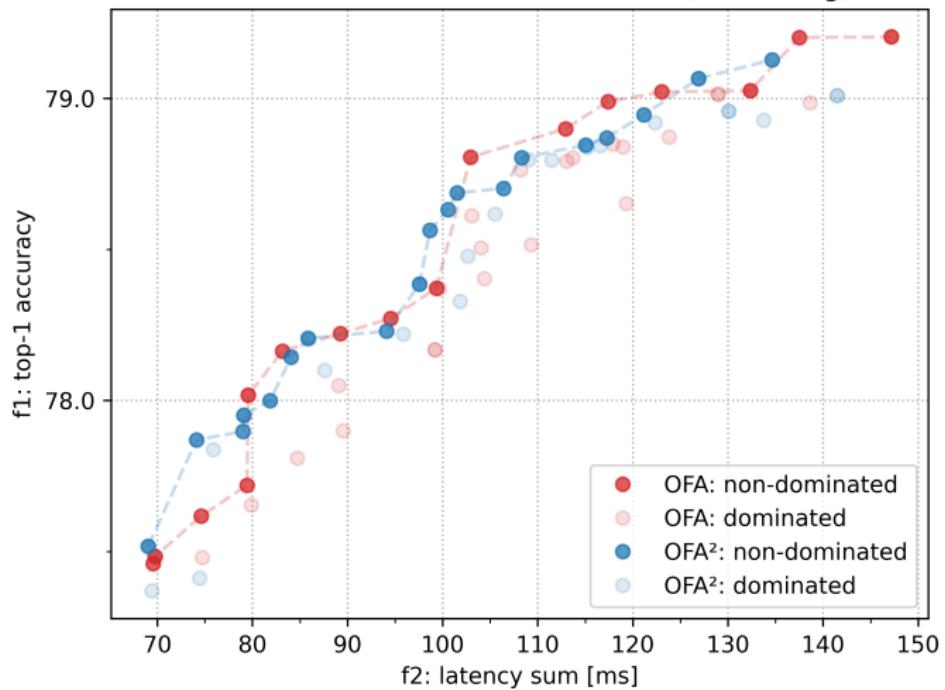
4 individuals

5 individuals

6 individuals

7 individuals

8 individuals

Ensemble OFA x OFA<sup>2</sup>: 3 individuals (soft voting)

# OFA<sup>2</sup> x OFA

- Voting

hard

soft

- Latency

sum

max

- Individuals

2 individuals

3 individuals

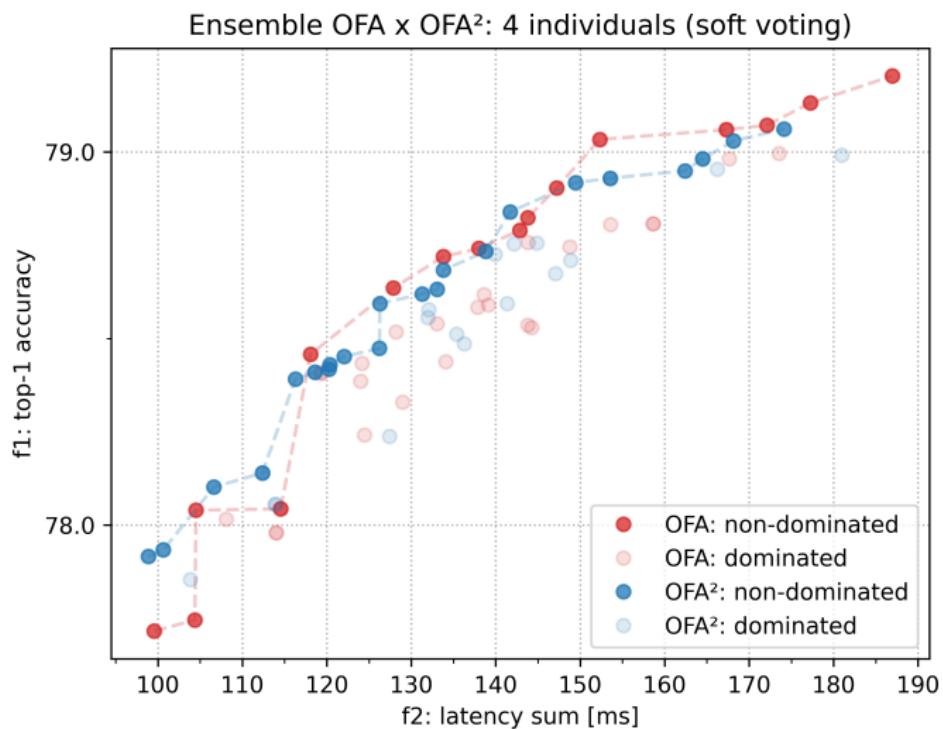
4 individuals

5 individuals

6 individuals

7 individuals

8 individuals



# OFA<sup>2</sup> x OFA

- Voting

hard

soft

- Latency

sum

max

- Individuals

2 individuals

3 individuals

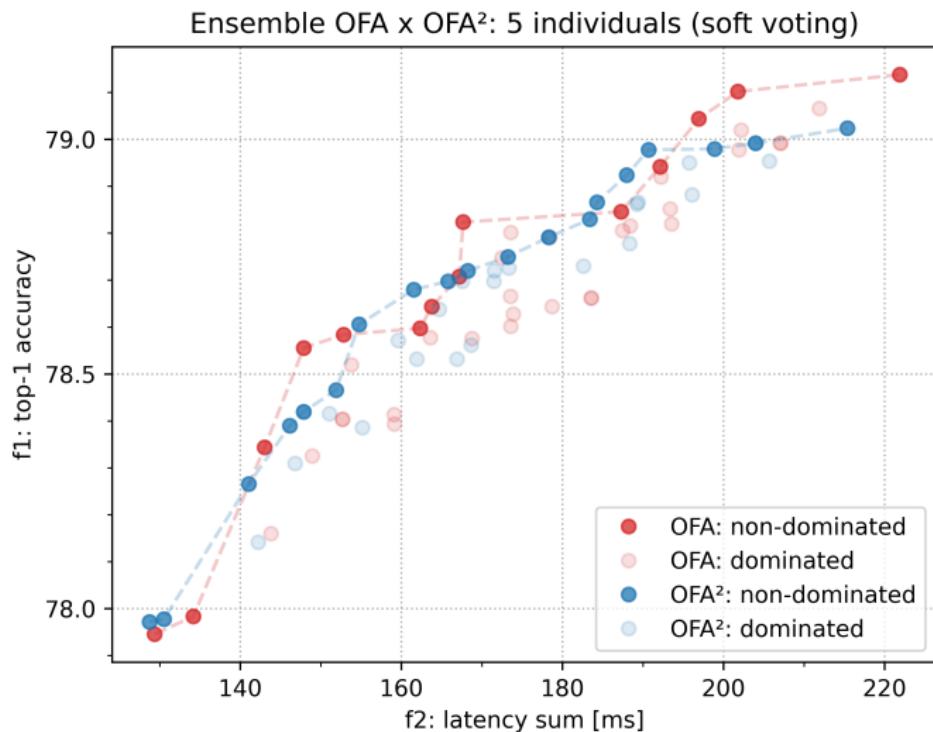
4 individuals

5 individuals

6 individuals

7 individuals

8 individuals



# OFA<sup>2</sup> x OFA

- Voting

hard

soft

- Latency

sum

max

- Individuals

2 individuals

3 individuals

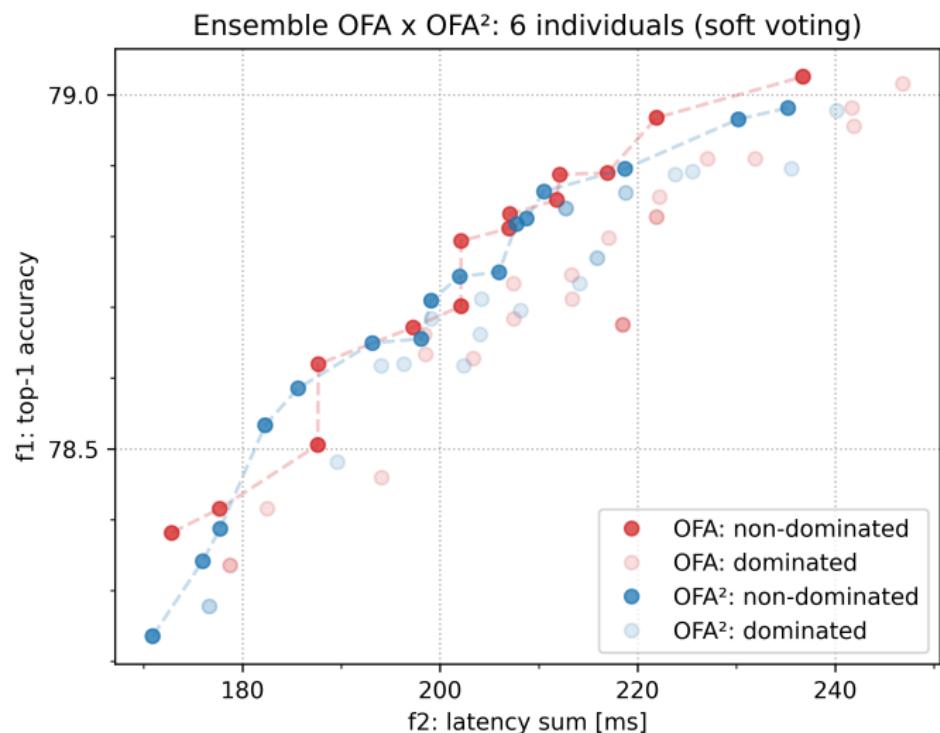
4 individuals

5 individuals

6 individuals

7 individuals

8 individuals



# OFA<sup>2</sup> × OFA

- Voting

hard

soft

- Latency

sum

max

- Individuals

2 individuals

3 individuals

4 individuals

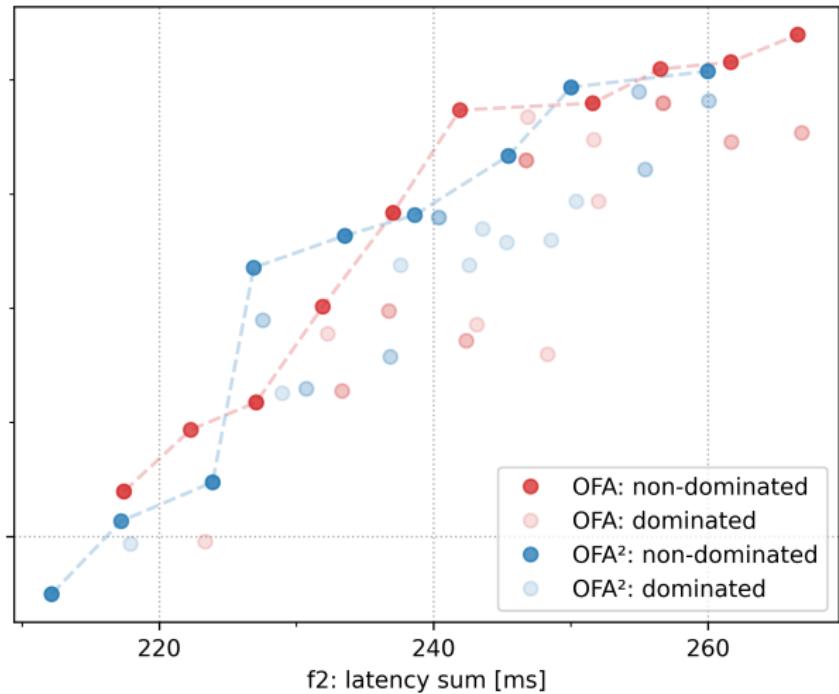
5 individuals

6 individuals

7 individuals

8 individuals

Ensemble OFA x OFA<sup>2</sup>: 7 individuals (soft voting)



# OFA<sup>2</sup> x OFA

- Voting

hard

soft

- Latency

sum

max

- Individuals

2 individuals

3 individuals

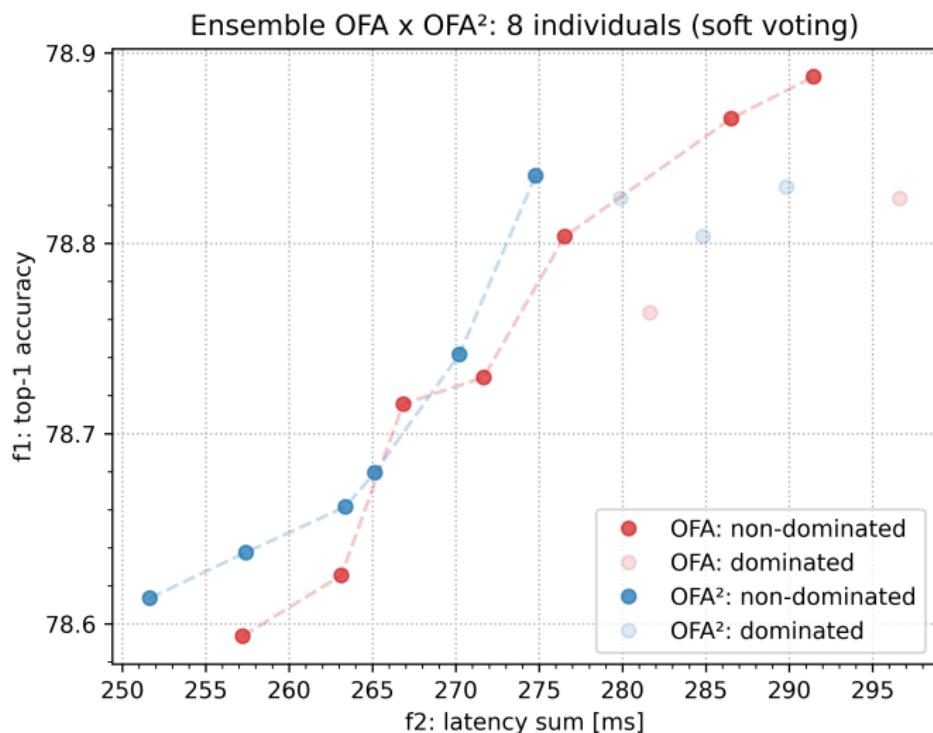
4 individuals

5 individuals

6 individuals

7 individuals

8 individuals



# OFA<sup>2</sup> × OFA

- **Voting**

hard

soft

- **Latency**

sum

max

- **Individuals**

2 individuals

3 individuals

4 individuals

5 individuals

6 individuals

7 individuals

8 individuals

# OFA<sup>2</sup> x OFA

- Voting

hard

soft

- Latency

sum

max

- Individuals

2 individuals

3 individuals

4 individuals

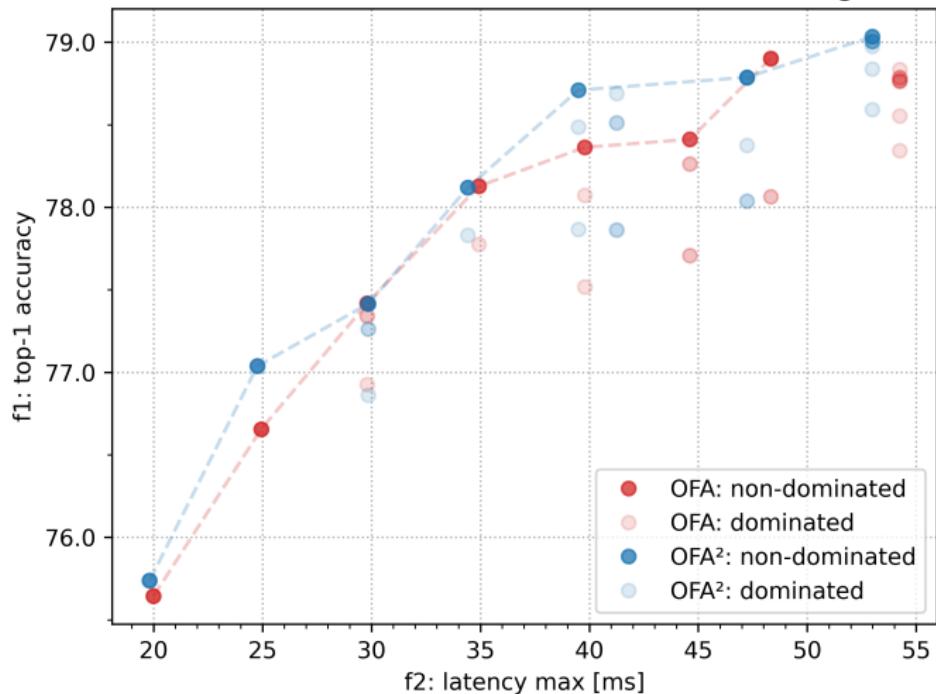
5 individuals

6 individuals

7 individuals

8 individuals

Ensemble OFA x OFA<sup>2</sup>: 2 individuals (hard voting)



# OFA<sup>2</sup> x OFA

- Voting

hard

soft

- Latency

sum

max

- Individuals

2 individuals

3 individuals

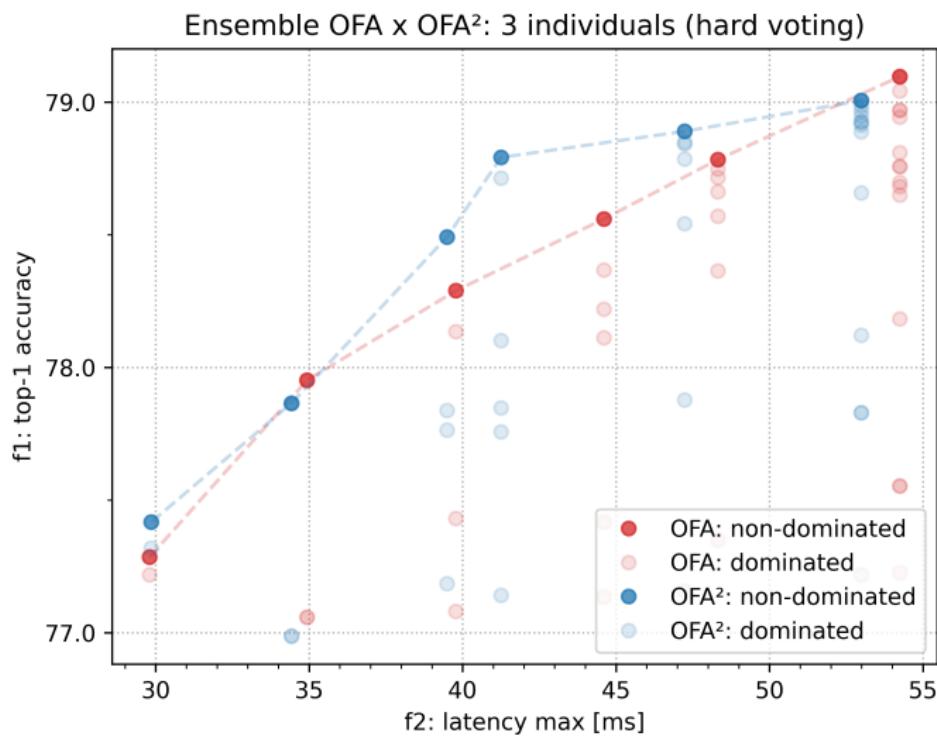
4 individuals

5 individuals

6 individuals

7 individuals

8 individuals



# OFA<sup>2</sup> x OFA

- Voting

hard

soft

- Latency

sum

max

- Individuals

2 individuals

3 individuals

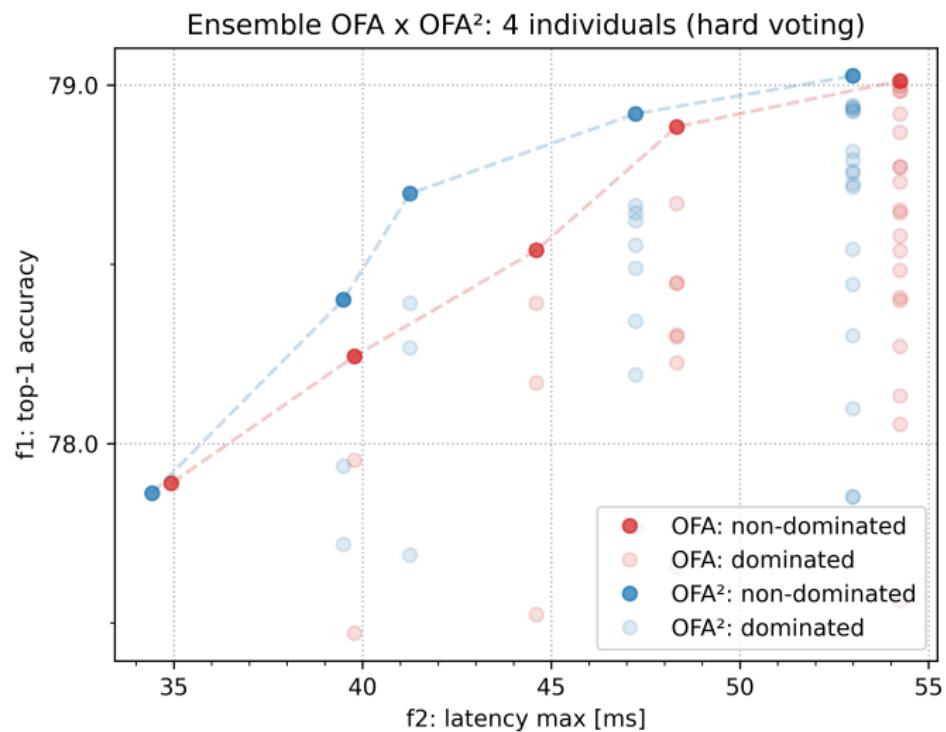
**4 individuals**

5 individuals

6 individuals

7 individuals

8 individuals



## OFA<sup>2</sup> × OFA

- Voting

hard

soft

- Latency

sum

max

- Individuals

2 individuals

### 3 individuals

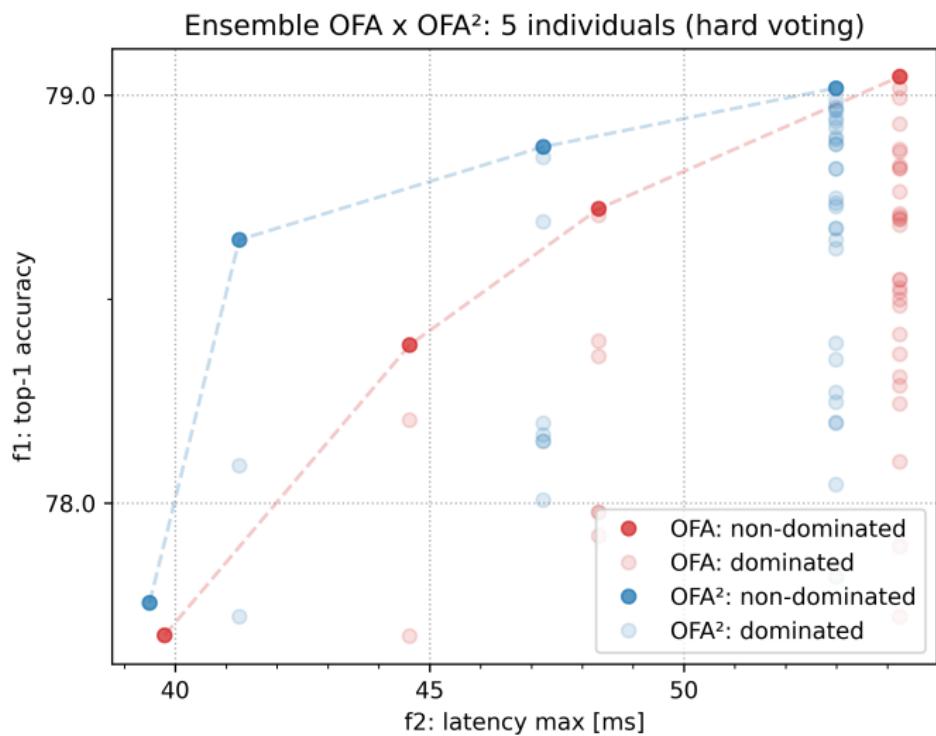
## 4 individuals

5 individuals

6 individuals

7 individuals

8 individuals



# OFA<sup>2</sup> x OFA

- Voting

hard

soft

- Latency

sum

max

- Individuals

2 individuals

3 individuals

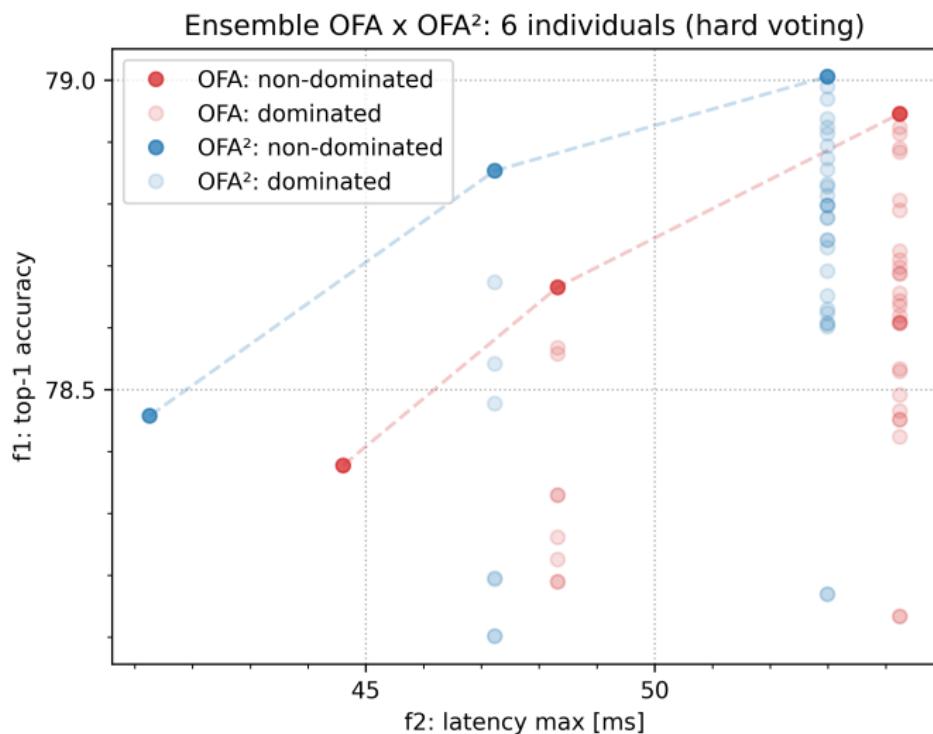
4 individuals

5 individuals

**6 individuals**

7 individuals

8 individuals



# OFA<sup>2</sup> x OFA

- Voting

hard

soft

- Latency

sum

max

- Individuals

2 individuals

3 individuals

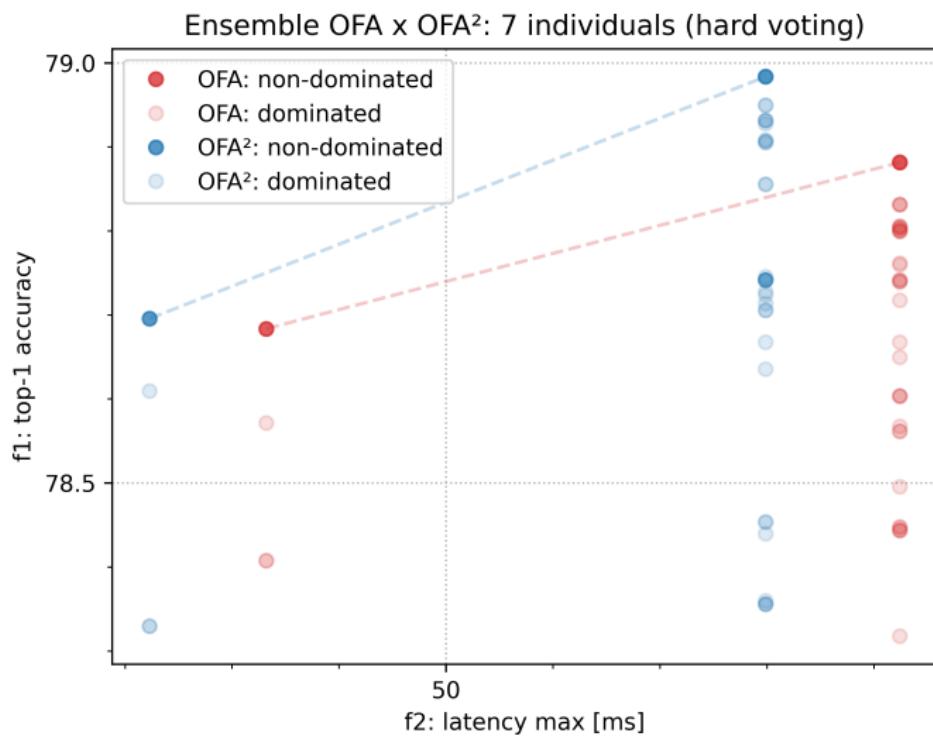
4 individuals

5 individuals

6 individuals

7 individuals

8 individuals



# OFA<sup>2</sup> x OFA

- Voting

hard

soft

- Latency

sum

max

- Individuals

2 individuals

3 individuals

4 individuals

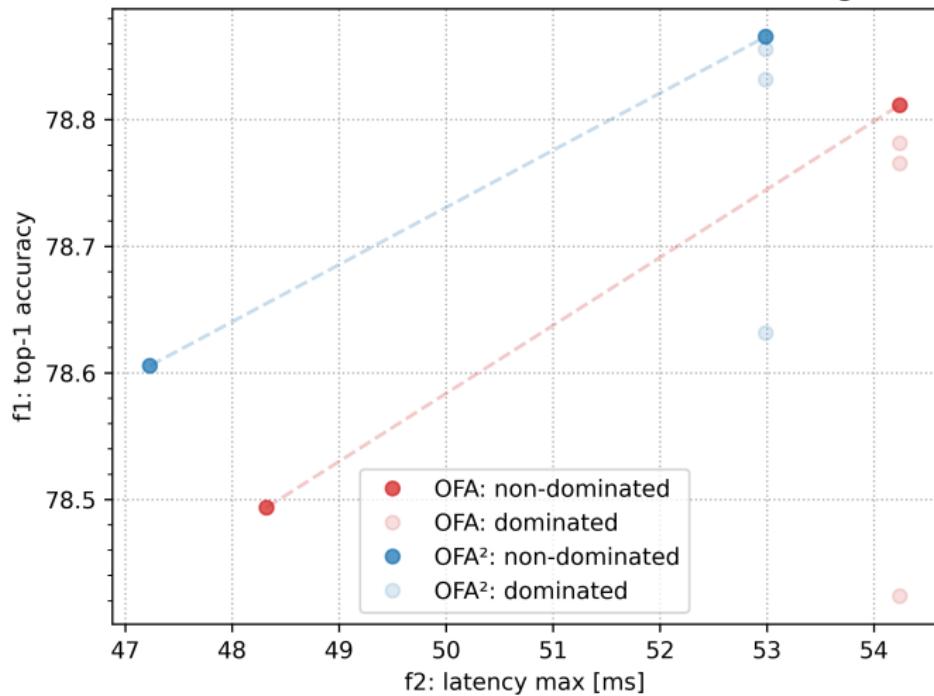
5 individuals

6 individuals

7 individuals

8 individuals

Ensemble OFA x OFA<sup>2</sup>: 8 individuals (hard voting)



# OFA<sup>2</sup> x OFA

- **Voting**

hard

soft

- **Latency**

sum

max

- **Individuals**

2 individuals

3 individuals

4 individuals

5 individuals

6 individuals

7 individuals

8 individuals

# OFA<sup>2</sup> x OFA

- Voting

hard

soft

- Latency

sum

max

- Individuals

2 individuals

3 individuals

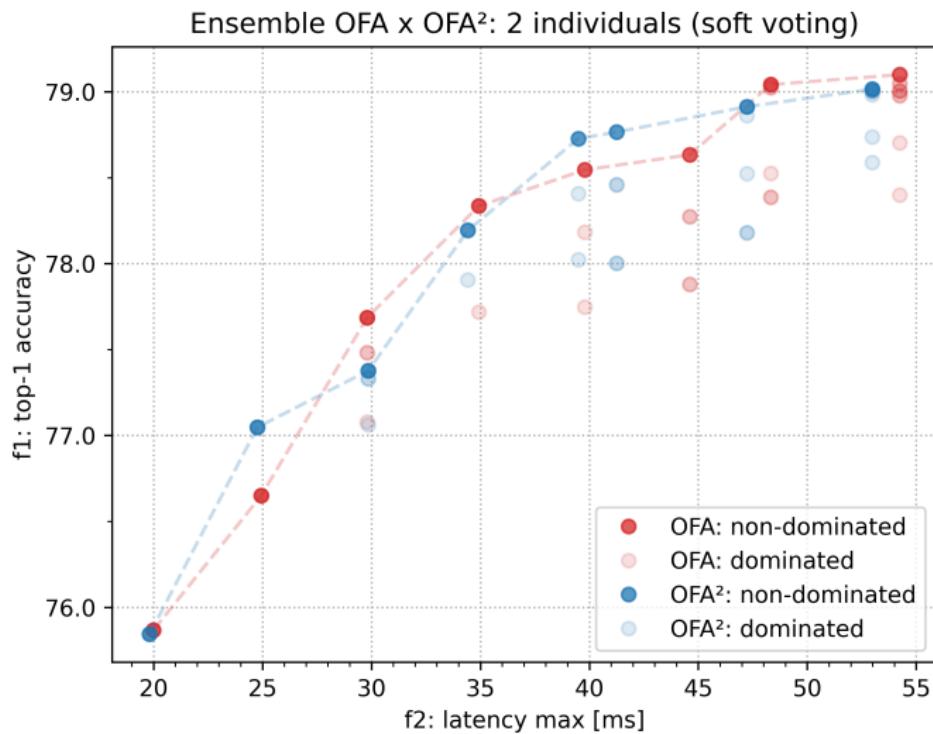
4 individuals

5 individuals

6 individuals

7 individuals

8 individuals



# OFA<sup>2</sup> x OFA

- Voting

hard

soft

- Latency

sum

max

- Individuals

2 individuals

3 individuals

4 individuals

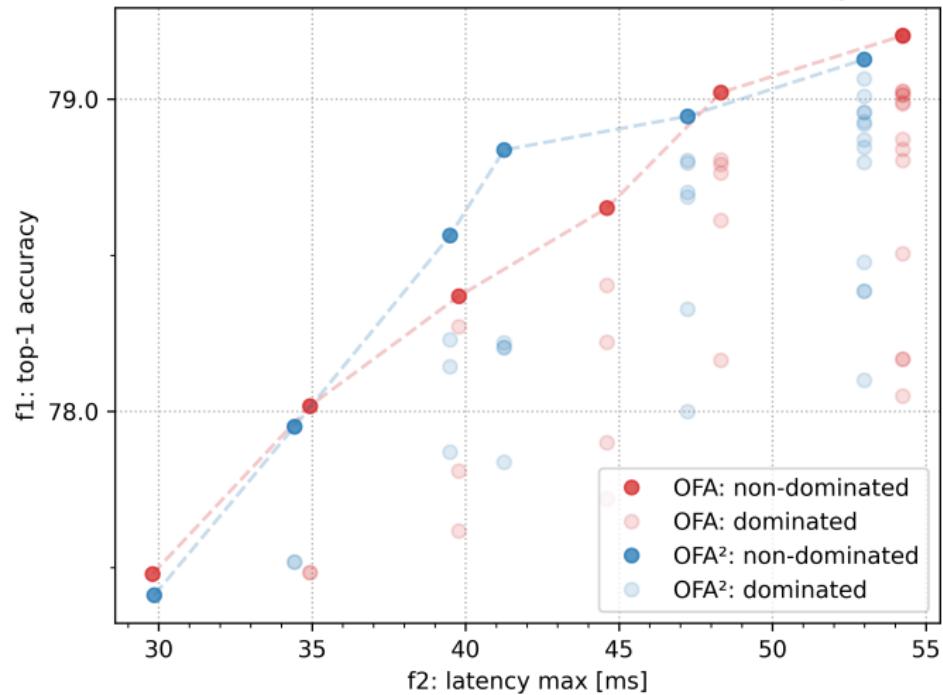
5 individuals

6 individuals

7 individuals

8 individuals

Ensemble OFA x OFA<sup>2</sup>: 3 individuals (soft voting)



# OFA<sup>2</sup> x OFA

- Voting

hard

soft

- Latency

sum

max

- Individuals

2 individuals

3 individuals

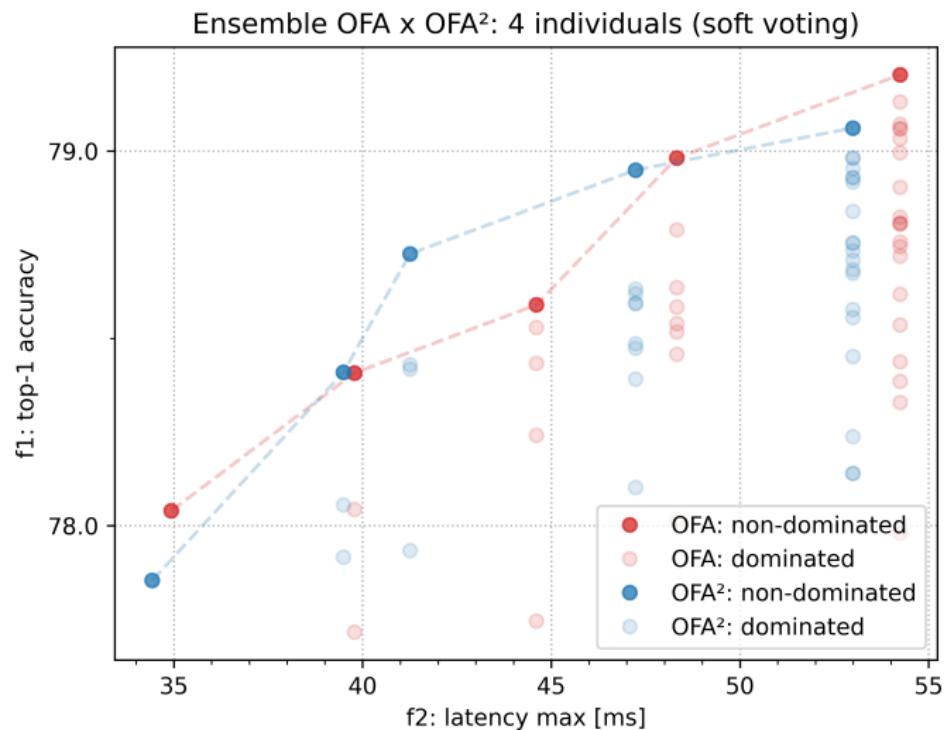
4 individuals

5 individuals

6 individuals

7 individuals

8 individuals



# OFA<sup>2</sup> x OFA

- Voting

hard

soft

- Latency

sum

max

- Individuals

2 individuals

3 individuals

4 individuals

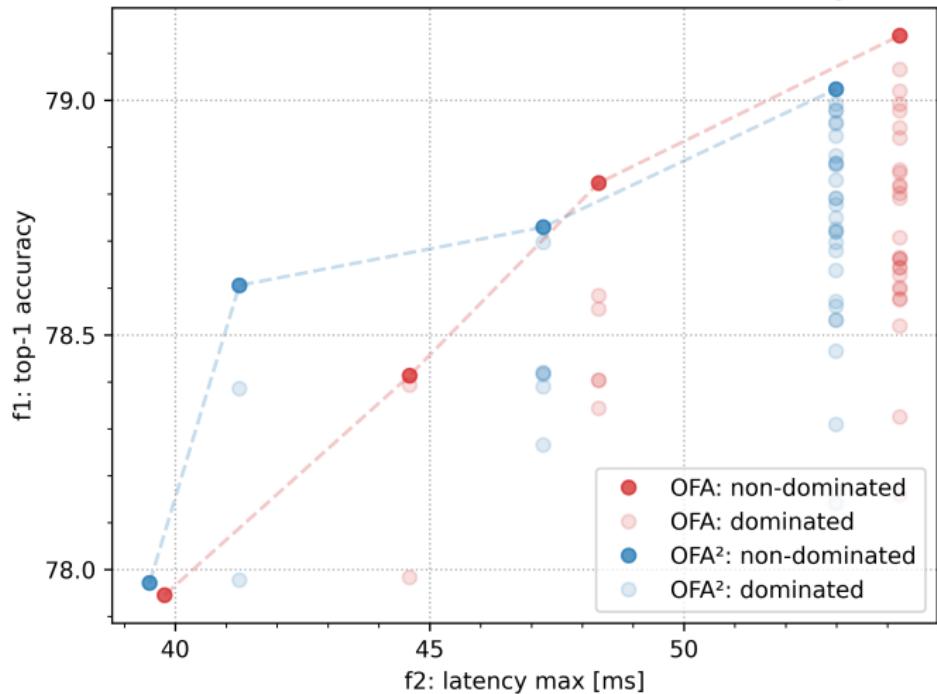
5 individuals

6 individuals

7 individuals

8 individuals

Ensemble OFA x OFA<sup>2</sup>: 5 individuals (soft voting)



# OFA<sup>2</sup> x OFA

- Voting

hard

soft

- Latency

sum

max

- Individuals

2 individuals

3 individuals

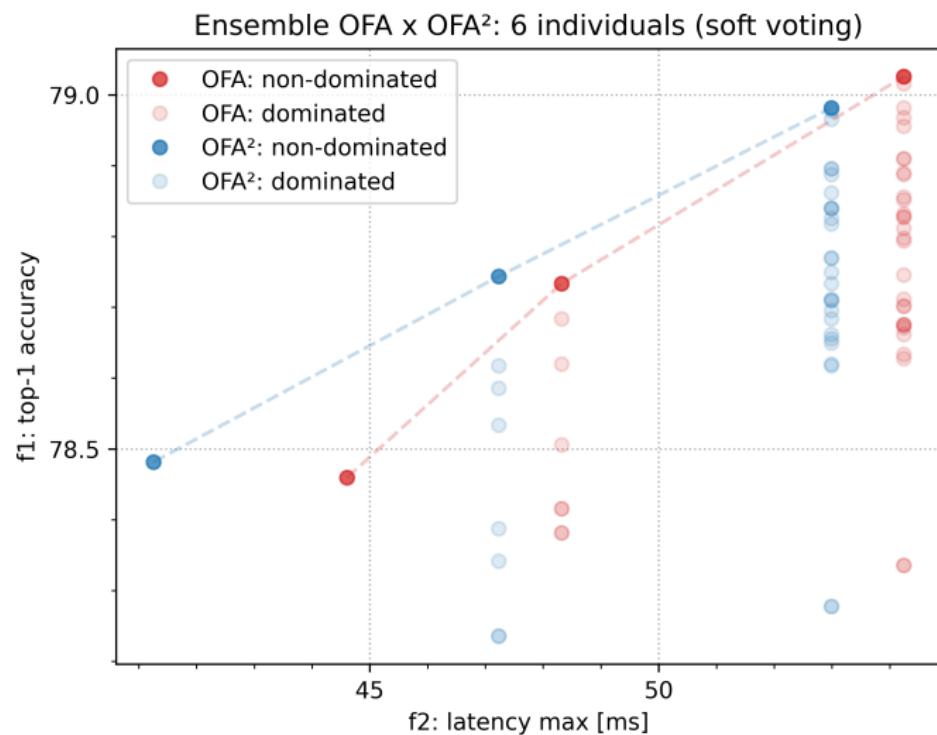
4 individuals

5 individuals

6 individuals

7 individuals

8 individuals



# OFA<sup>2</sup> x OFA

- Voting

hard

soft

- Latency

sum

max

- Individuals

2 individuals

3 individuals

4 individuals

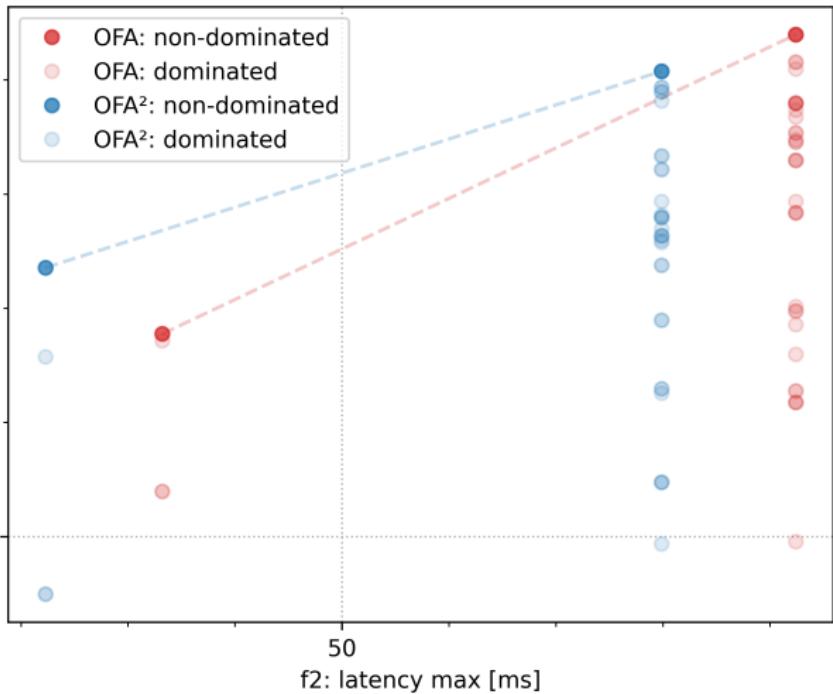
5 individuals

6 individuals

7 individuals

8 individuals

Ensemble OFA x OFA<sup>2</sup>: 7 individuals (soft voting)



# OFA<sup>2</sup> x OFA

- Voting

hard

soft

- Latency

sum

max

- Individuals

2 individuals

3 individuals

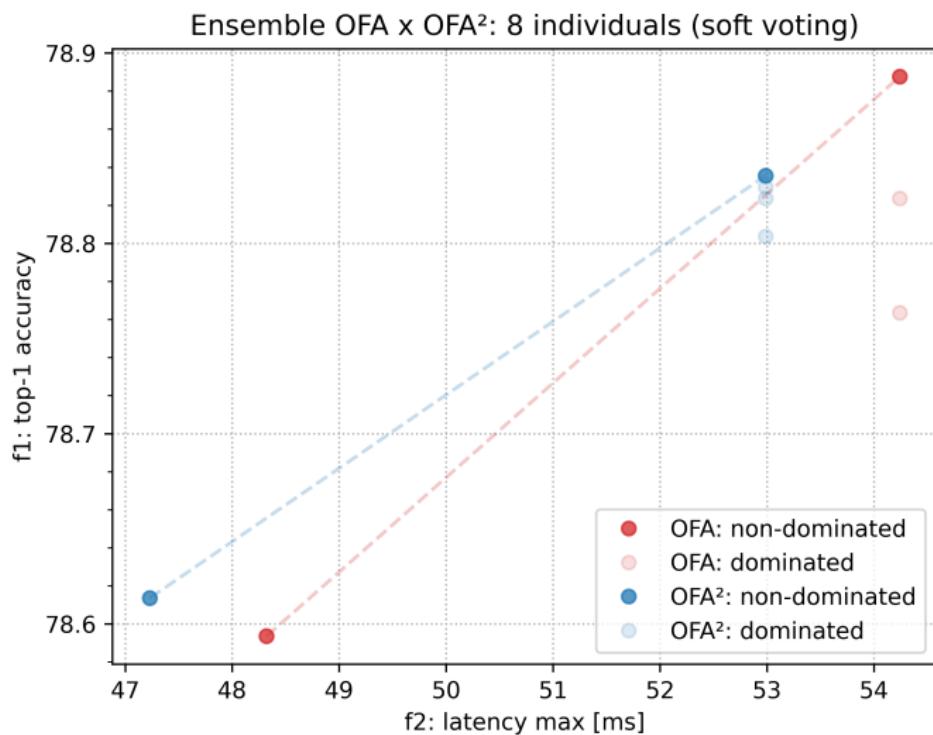
4 individuals

5 individuals

6 individuals

7 individuals

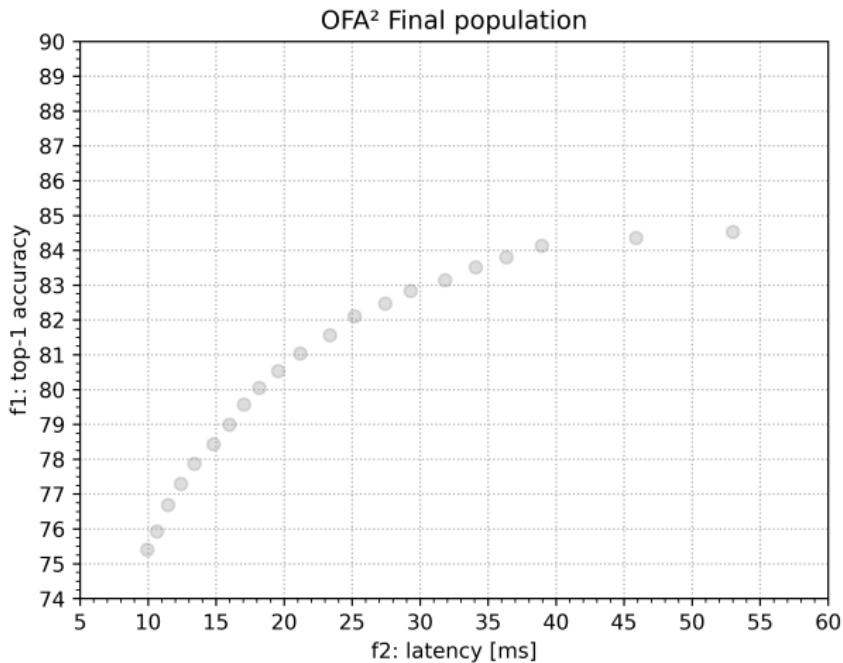
8 individuals



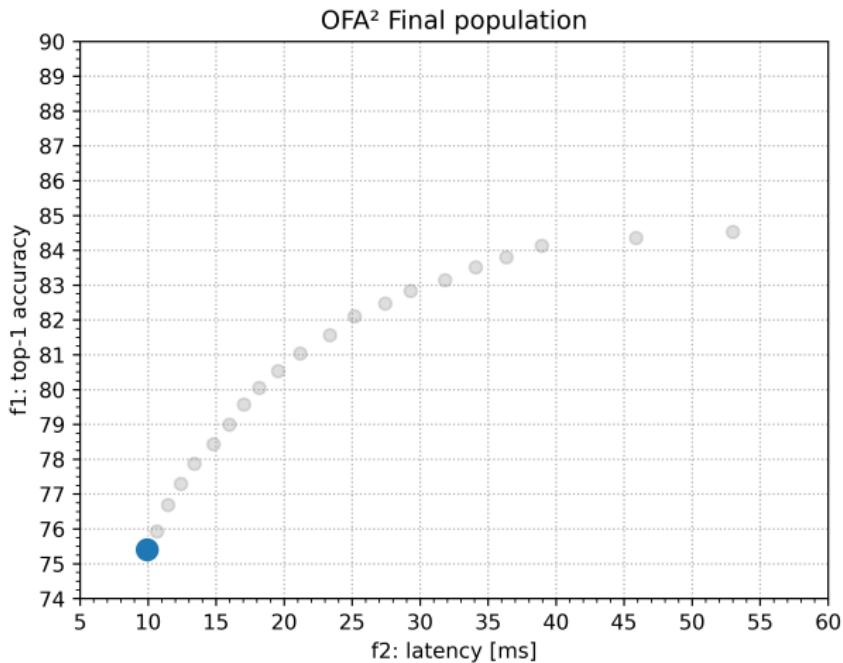
# Agenda

- 1 Introduction
- 2 AutoML
- 3 Training stage
- 4 Search stage
  - Random search
  - OFA search
  - OFA<sup>2</sup> search
  - Comparison
- 5 Ensemble
  - Manual sampling
  - Multi-objective Optimization
- 6 Conclusion
- 7 References

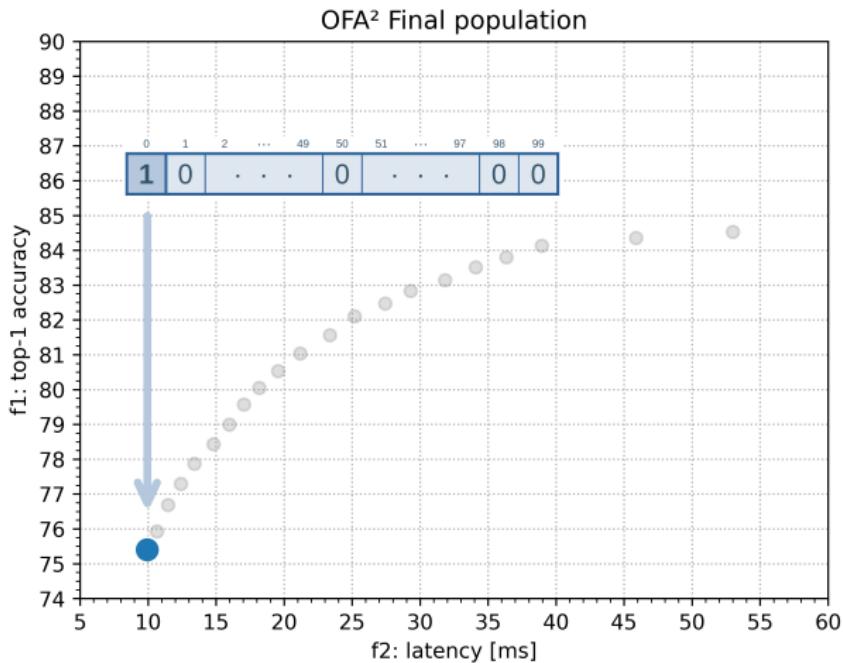
# Encoding



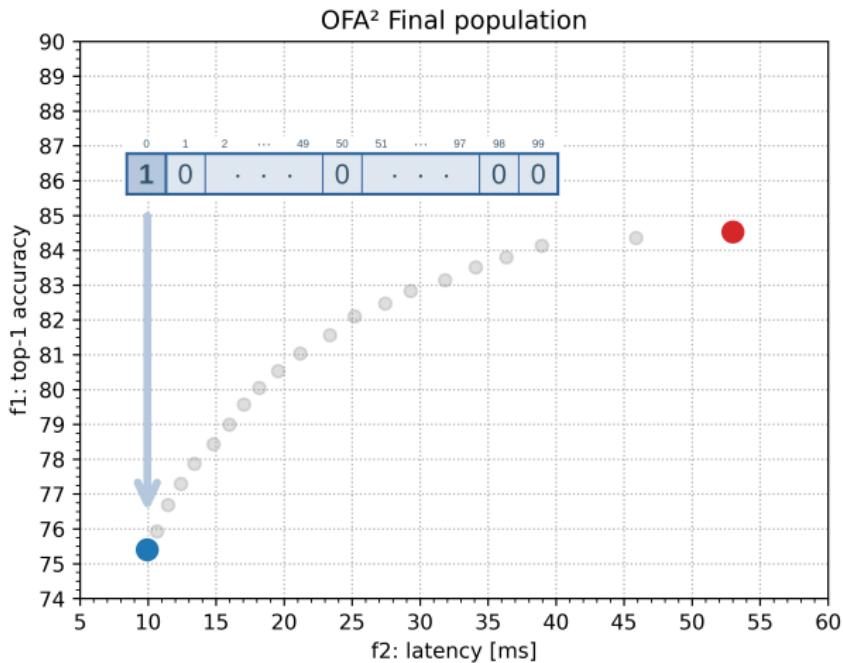
# Encoding



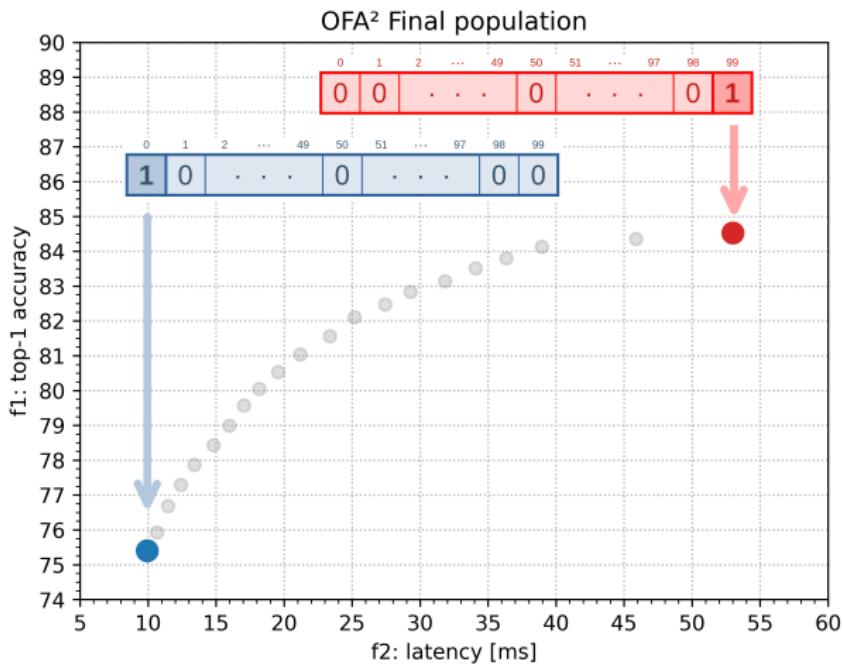
# Encoding



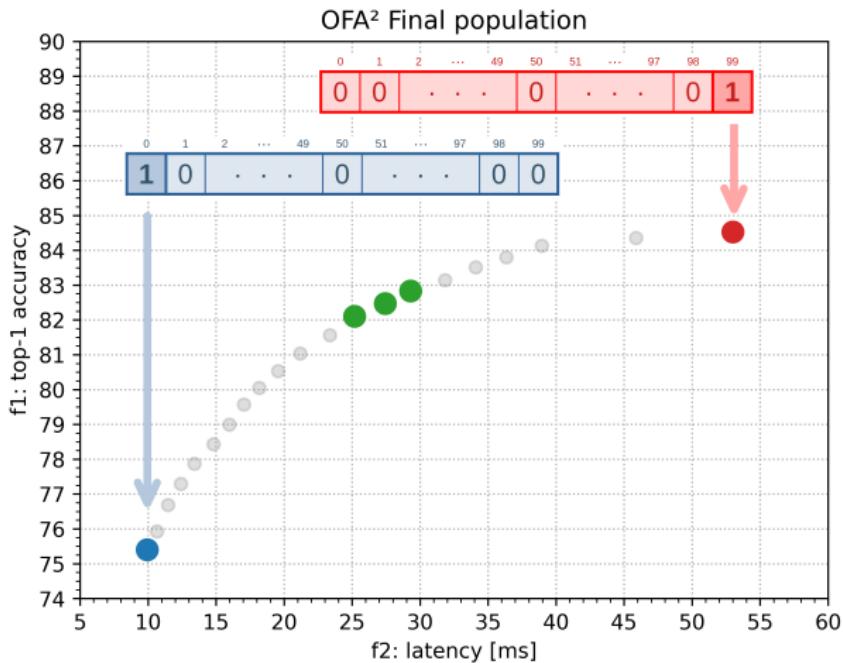
# Encoding



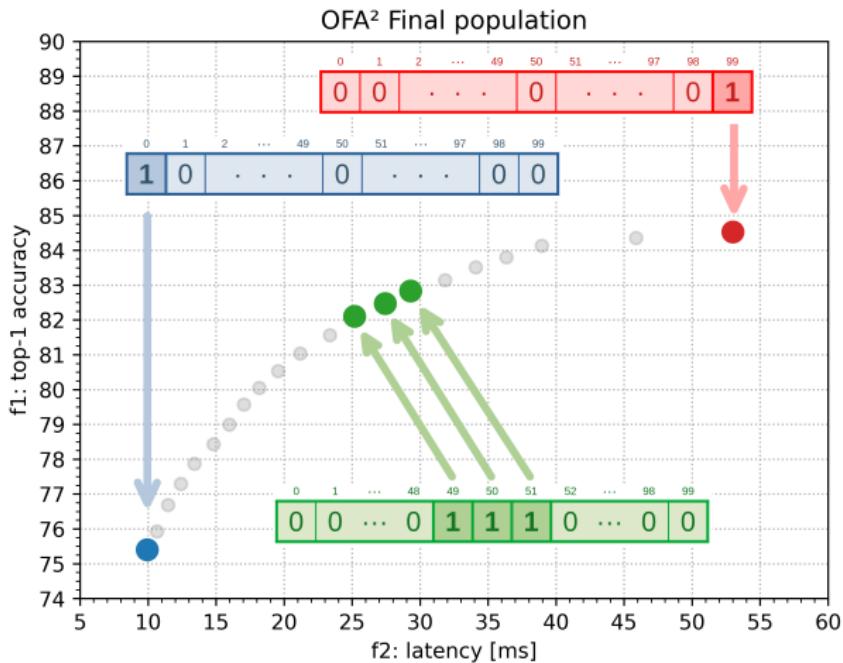
# Encoding



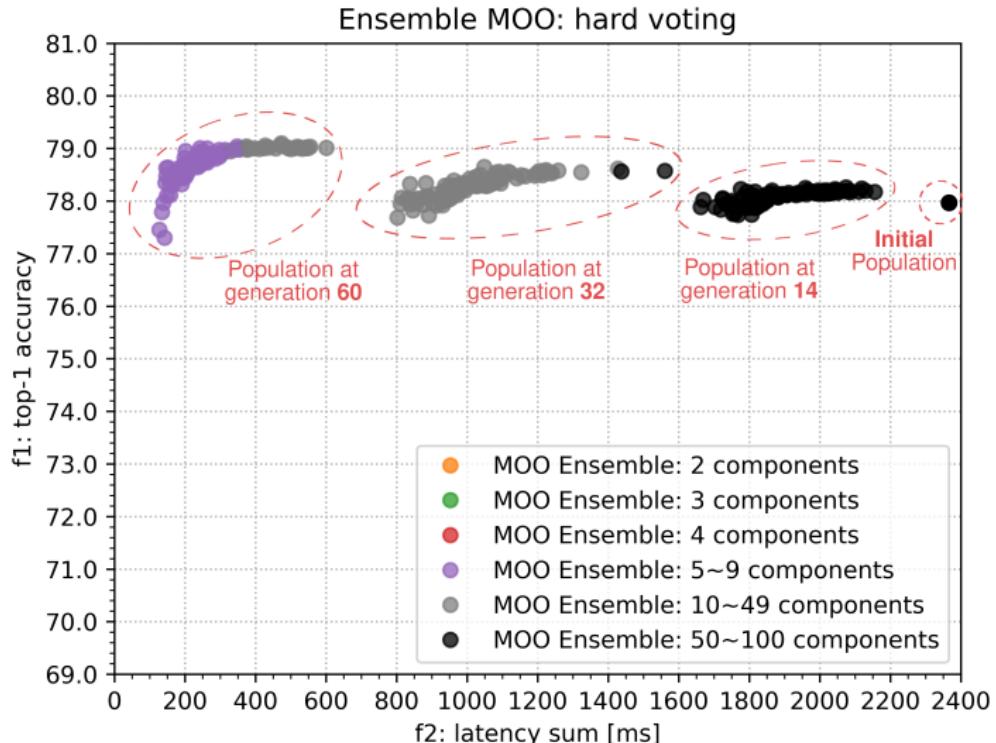
# Encoding



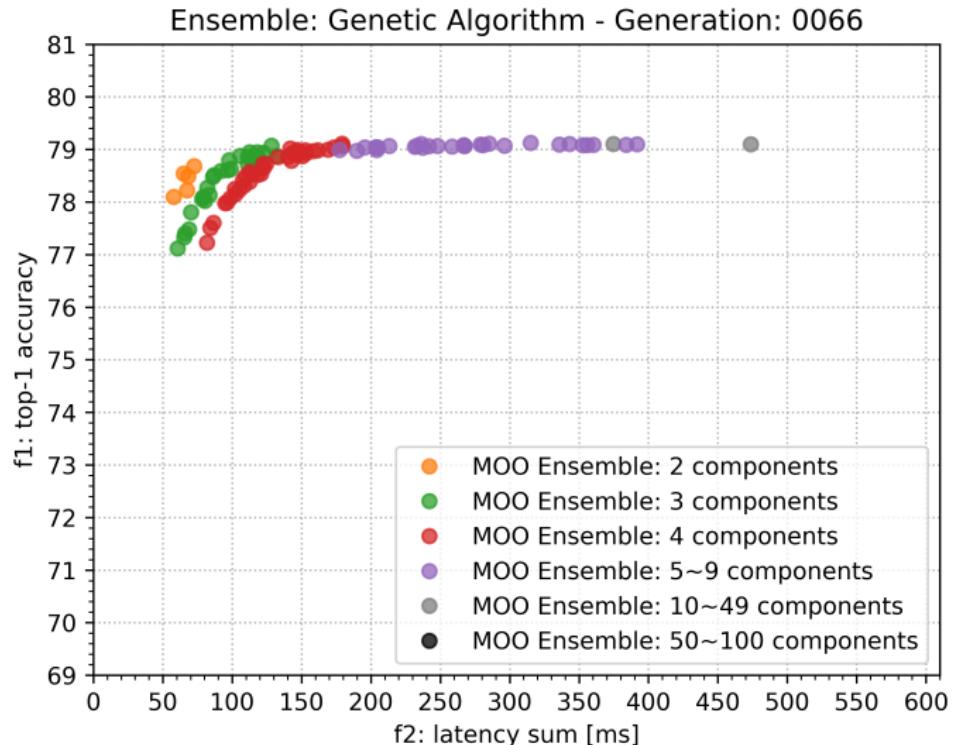
# Encoding



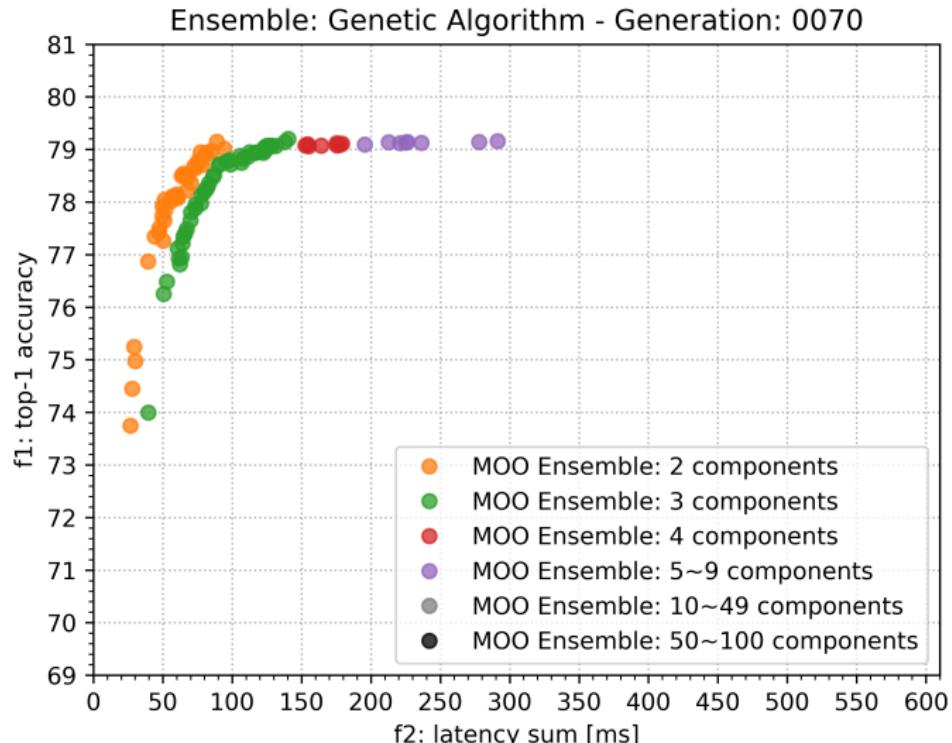
# Sum latency: Population progression



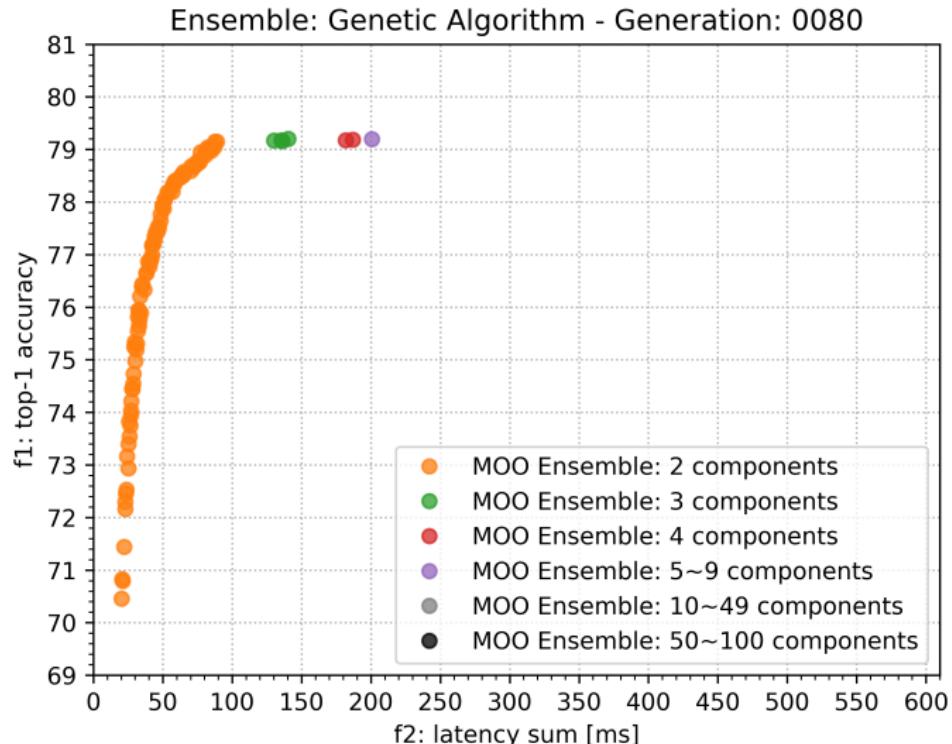
# Sum latency: Population progression



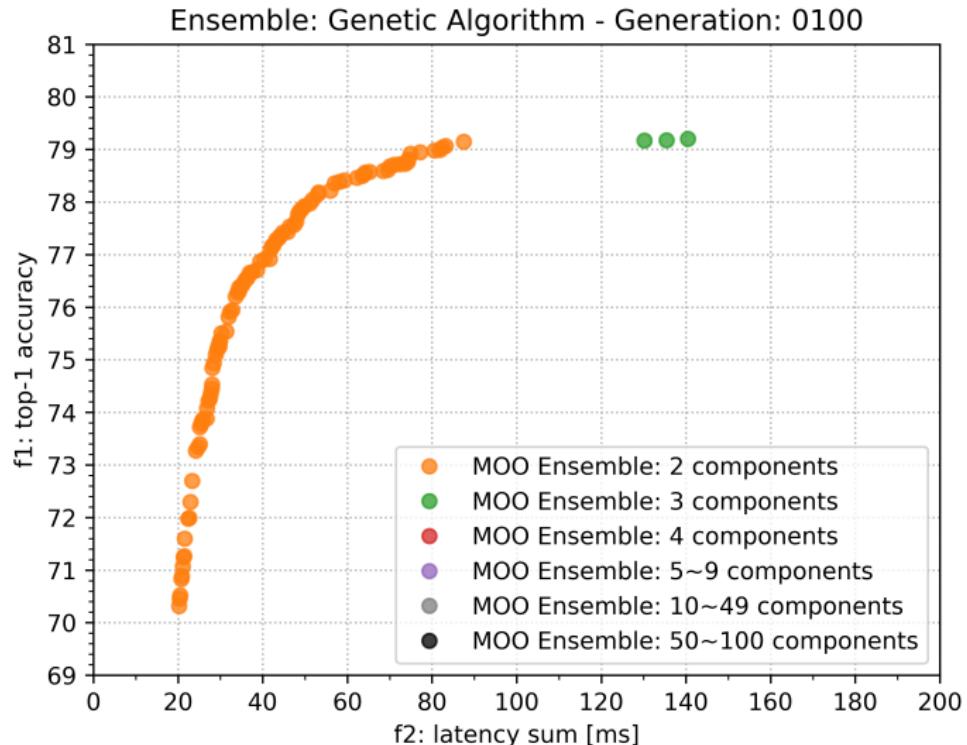
# Sum latency: Population progression



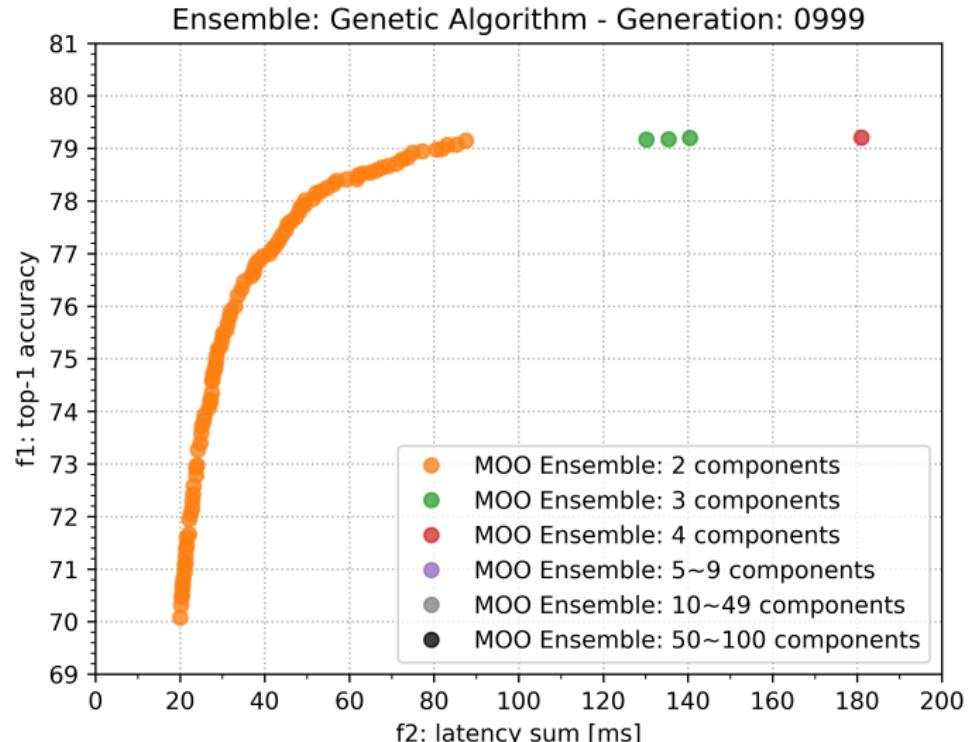
# Sum latency: Population progression



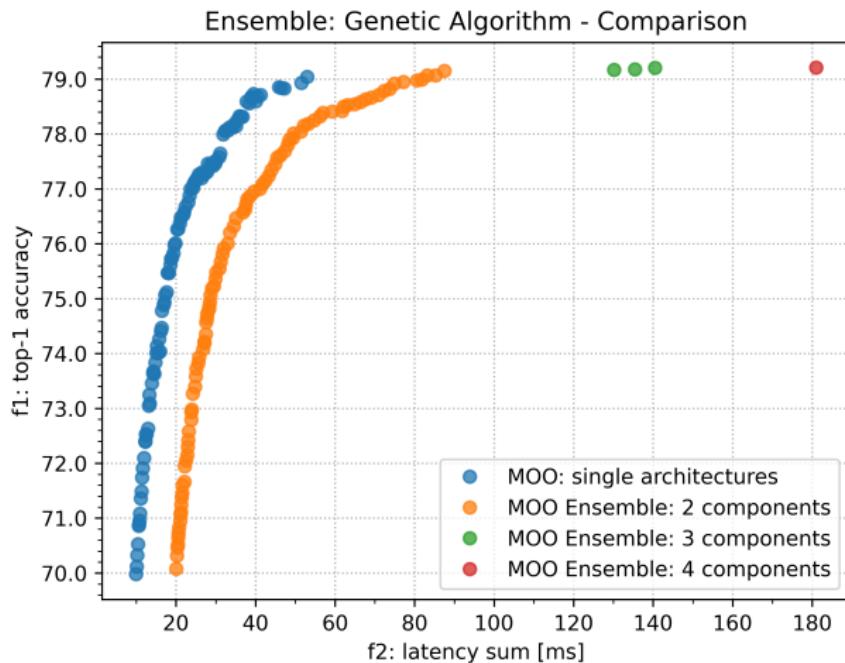
# Sum latency: Population progression



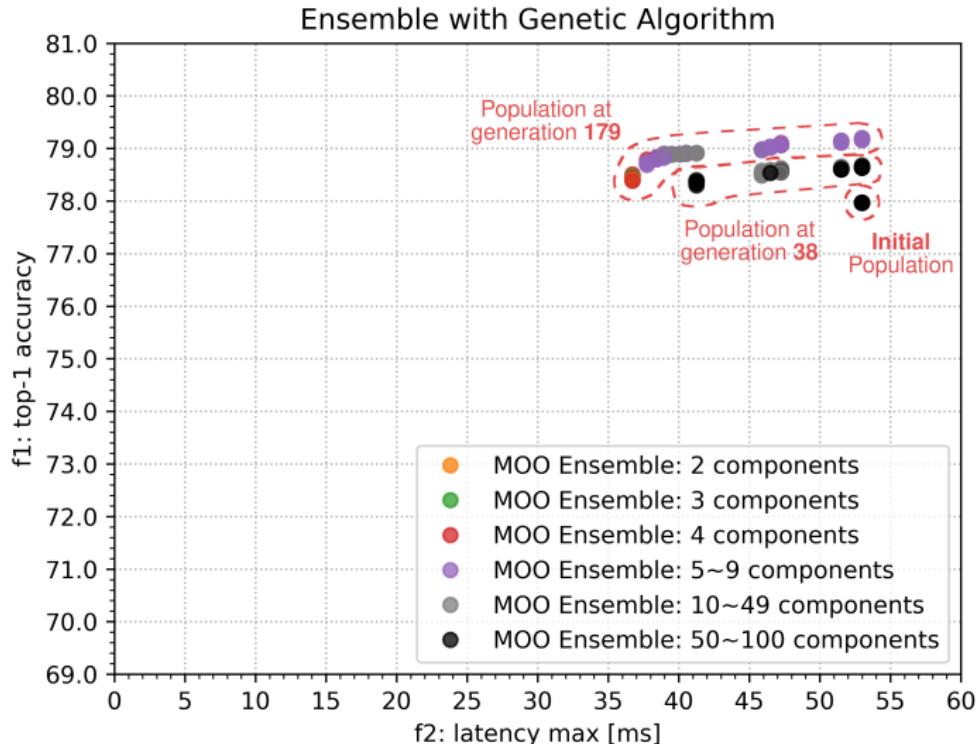
# Sum latency: Population progression



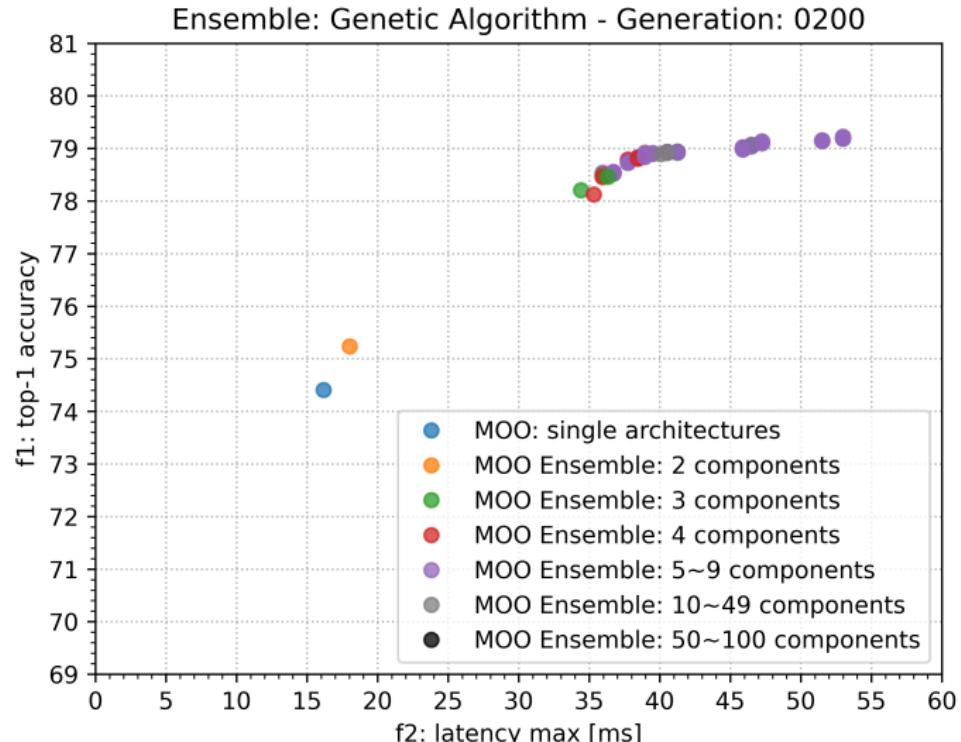
# Sum latency: Final population



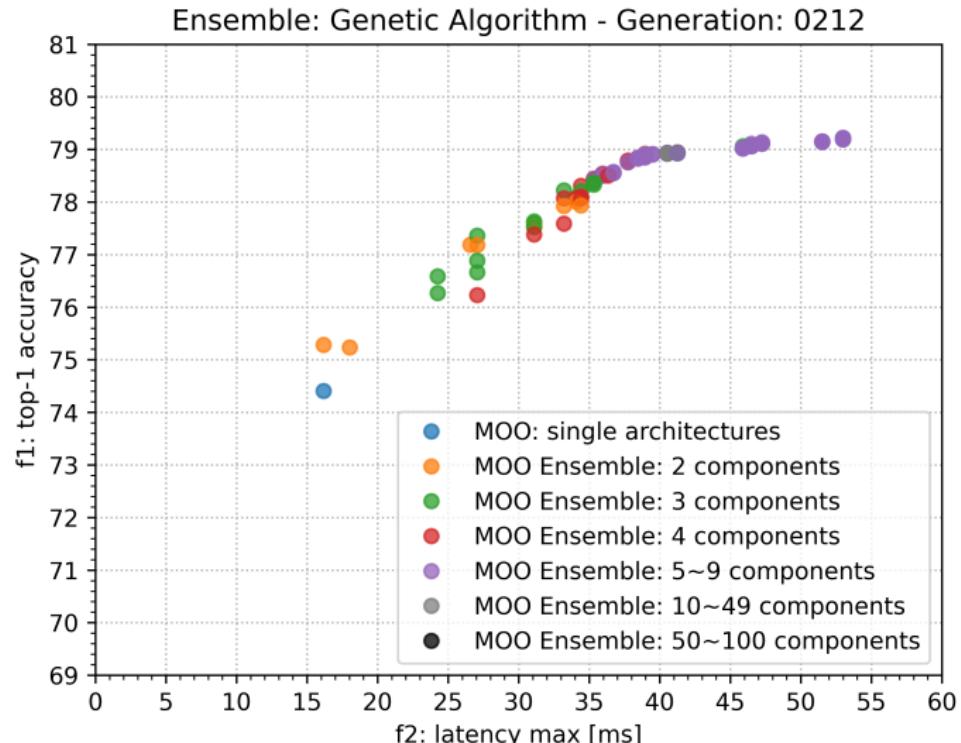
# Max latency: Population progression



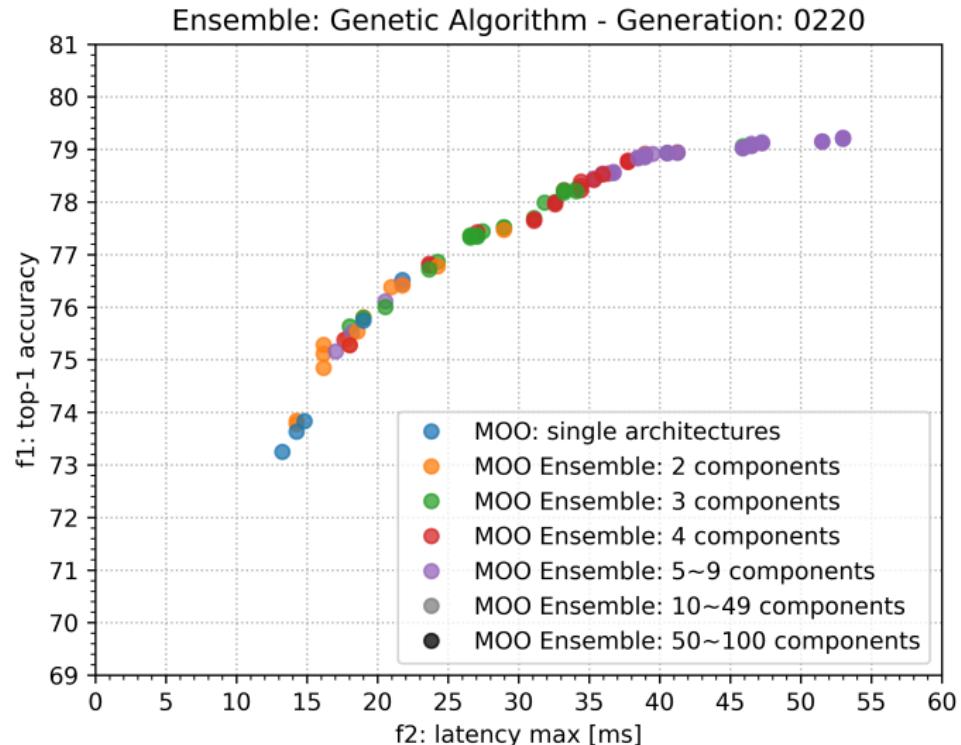
# Max latency: Population progression



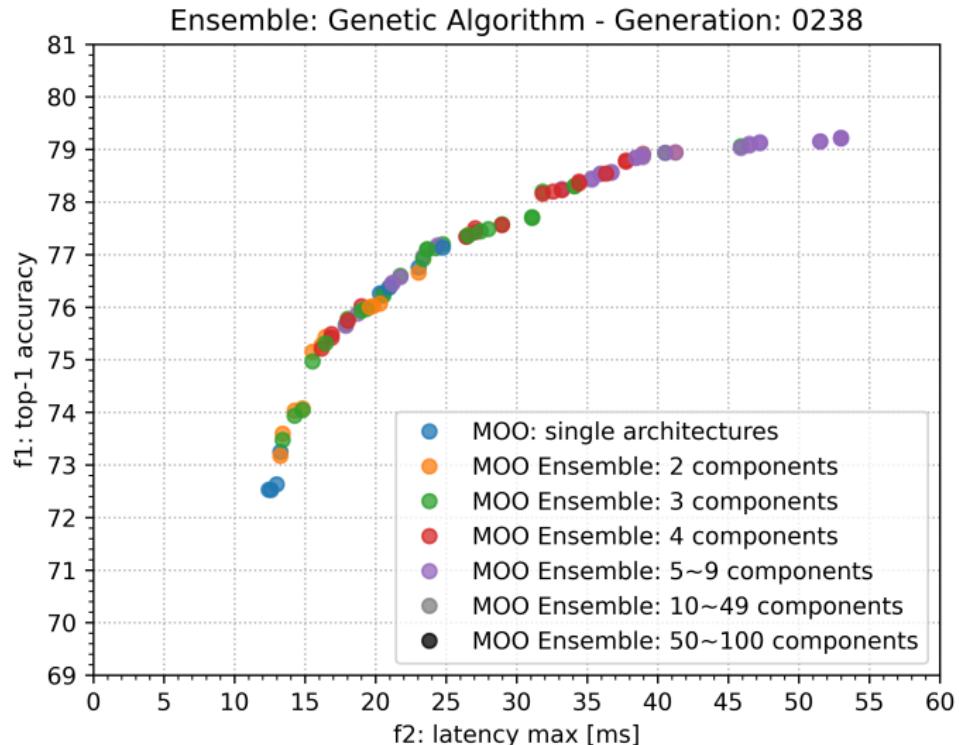
# Max latency: Population progression



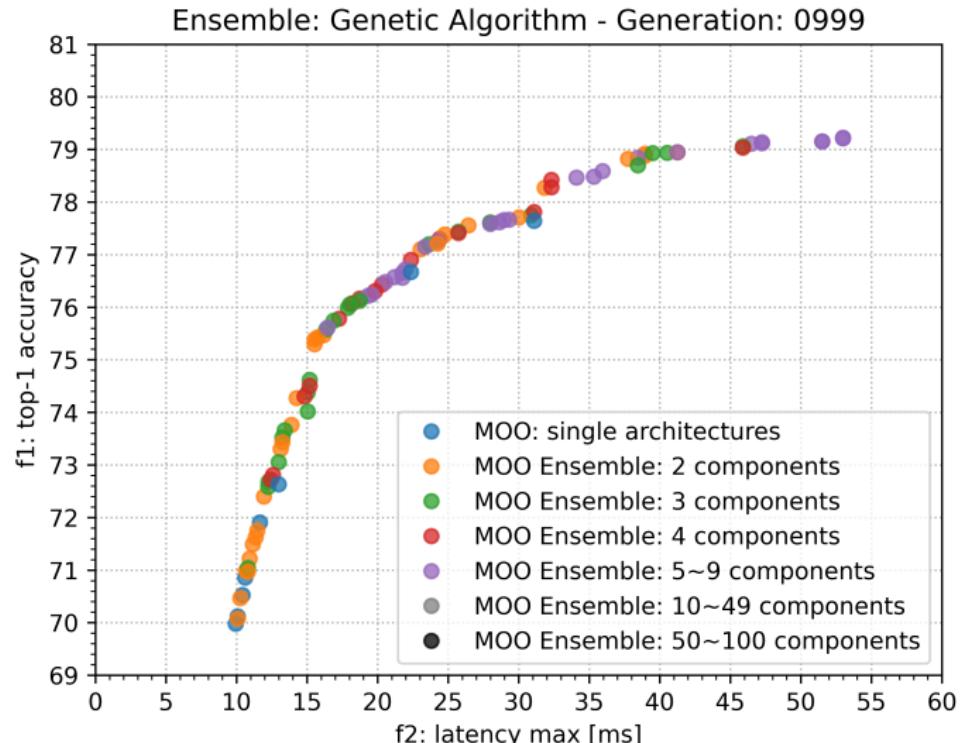
# Max latency: Population progression



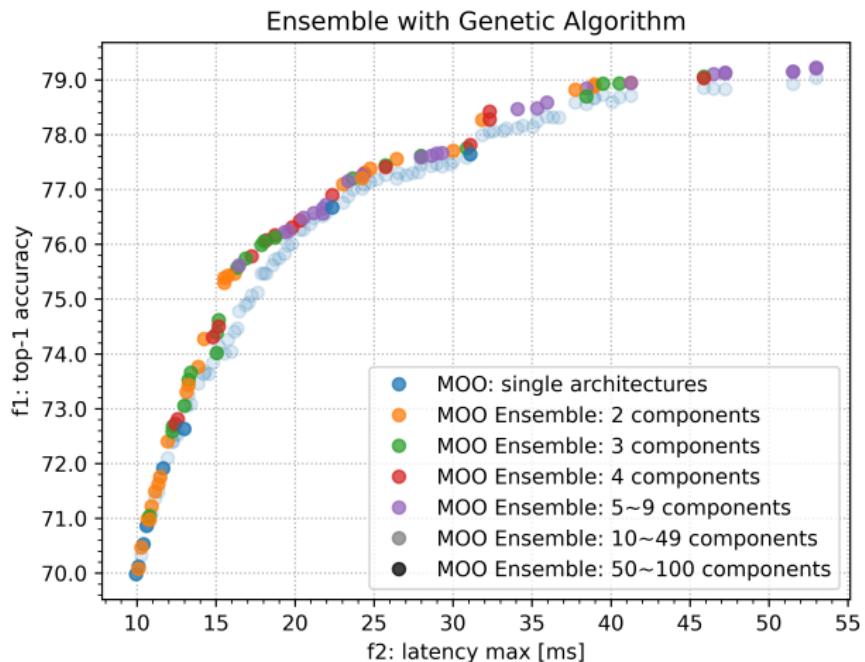
# Max latency: Population progression



# Max latency: Population progression



## Max latency: Final population



# Agenda

- 1 Introduction
- 2 AutoML
- 3 Training stage
- 4 Search stage
  - Random search
  - OFA search
  - OFA<sup>2</sup> search
  - Comparison
- 5 Ensemble
  - Manual sampling
  - Multi-objective Optimization
- 6 Conclusion
- 7 References

## Future work

- Fine-tuning the search using the real evaluated accuracies of the final population of neural networks obtained with the accuracy predictor.
- Perform the MOO search with all other hardware available.
- Generate accuracy predictor & latency lookup table for more different hardwares

## Summary

- OFA<sup>2</sup> improves the OFA framework by:
  - Finding better architectures with respect to the objectives functions.
  - Find a set of solutions all at once in a single search.
  - Keeps the computational cost roughly the same.
- Proposal of an EMOA encoding for the ensemble formation.
- No clear winner between hard/soft voting strategies for ensembles.
- Committee machine with better accuracies considering the maximum latency (parallel).
- Source code and experiments publicly available.

# Agenda

- 1 Introduction
- 2 AutoML
- 3 Training stage
- 4 Search stage
  - Random search
  - OFA search
  - OFA<sup>2</sup> search
  - Comparison
- 5 Ensemble
  - Manual sampling
  - Multi-objective Optimization
- 6 Conclusion
- 7 References

- Cai, Han, Chuang Gan, Tianzhe Wang, Zhekai Zhang, and Song Han. 2020. “Once for All: Train One Network and Specialize It for Efficient Deployment.” [https://iclr.cc/virtual\\_2020/poster\\_HylxE1HKwS.html](https://iclr.cc/virtual_2020/poster_HylxE1HKwS.html).
- Cai, Han, Ligeng Zhu, and Song Han. 2022. “ProxylessNAS: Direct Neural Architecture Search on Target Task and Hardware.” <https://openreview.net/forum?id=HylVB3AqYm>.
- Elsken, Thomas, Jan Hendrik Metzen, and Frank Hutter. 2019. “Neural Architecture Search: A Survey.” *Journal of Machine Learning Research* 20 (55): 1–21. <http://jmlr.org/papers/v20/18-598.html>.
- Green, Sam, Craig M. Vineyard, Ryan Helinski, and Çetin Kaya Koç. 2019. “RAPDARTS: Resource-Aware Progressive Differentiable Architecture Search.” <https://doi.org/10.48550/ARXIV.1911.05704>.
- Real, Esteban, Sherry Moore, Andrew Selle, Saurabh Saxena, Yutaka Leon Suematsu, Jie Tan, Quoc V. Le, and Alexey Kurakin. 2017. “Large-Scale Evolution of Image Classifiers.” In *Proceedings of the 34th International Conference on Machine Learning*, 2902–11. PMLR. <https://proceedings.mlr.press/v70/real17a.html>.

- Ying, Chris, Aaron Klein, Eric Christiansen, Esteban Real, Kevin Murphy, and Frank Hutter. 2019. "NAS-Bench-101: Towards Reproducible Neural Architecture Search." In *Proceedings of the 36th International Conference on Machine Learning*, 7105–14. PMLR.  
<https://proceedings.mlr.press/v97/ying19a.html>.