# Time series classification at scale
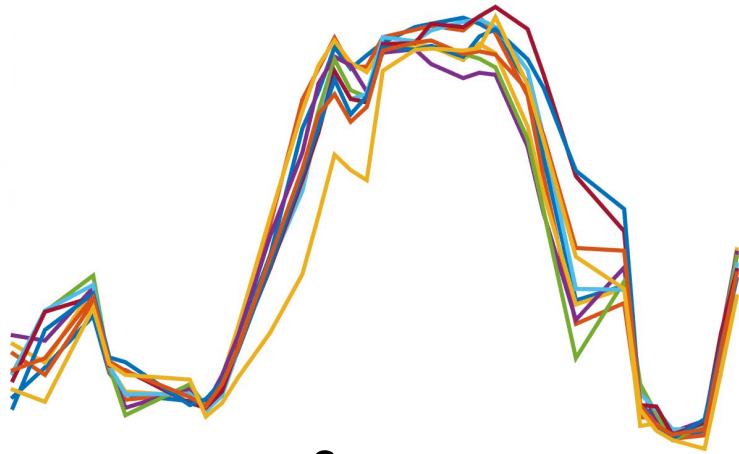
## François Petitjean

http://francois-petitjean.com

PhD work of A. Shifaz, B. Lucas and H. Ismail Fawaz
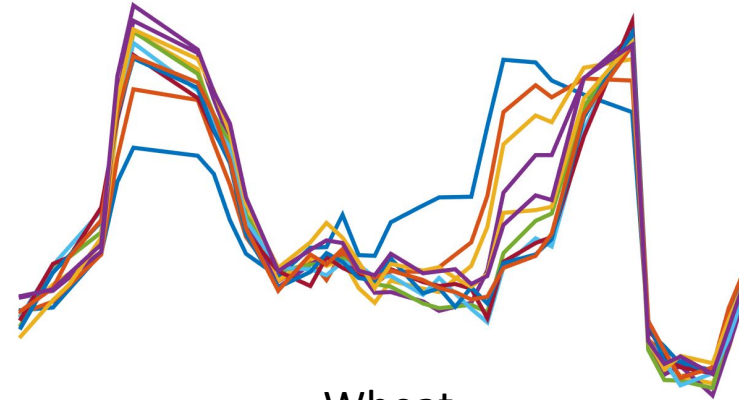
… with G. Forestier, C. Pelletier, G. Webb, J. Weber, B. Goethals, J. Weber, D. Schmidt, L. Idoumghar, P. Muller, L. O'Neill, and N. Zaidi

MONASH University
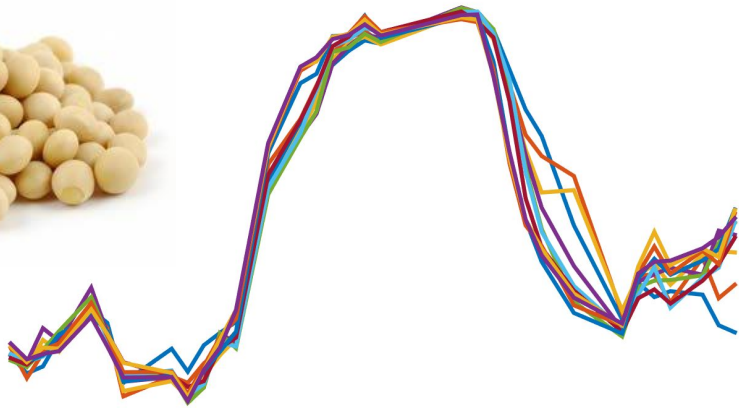
GROUP OF EIGHT AUSTRALIA

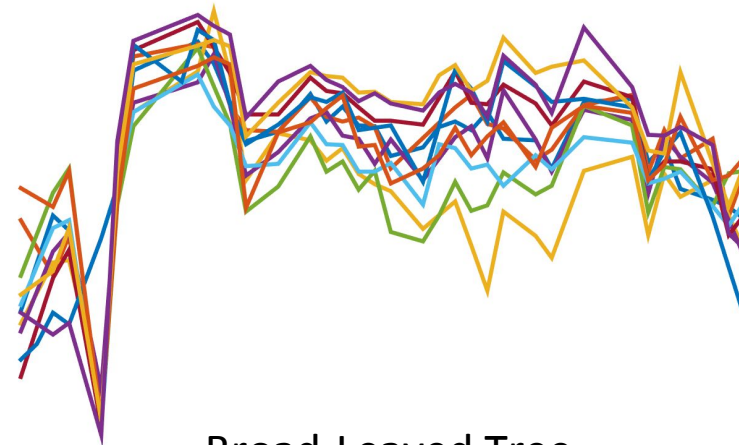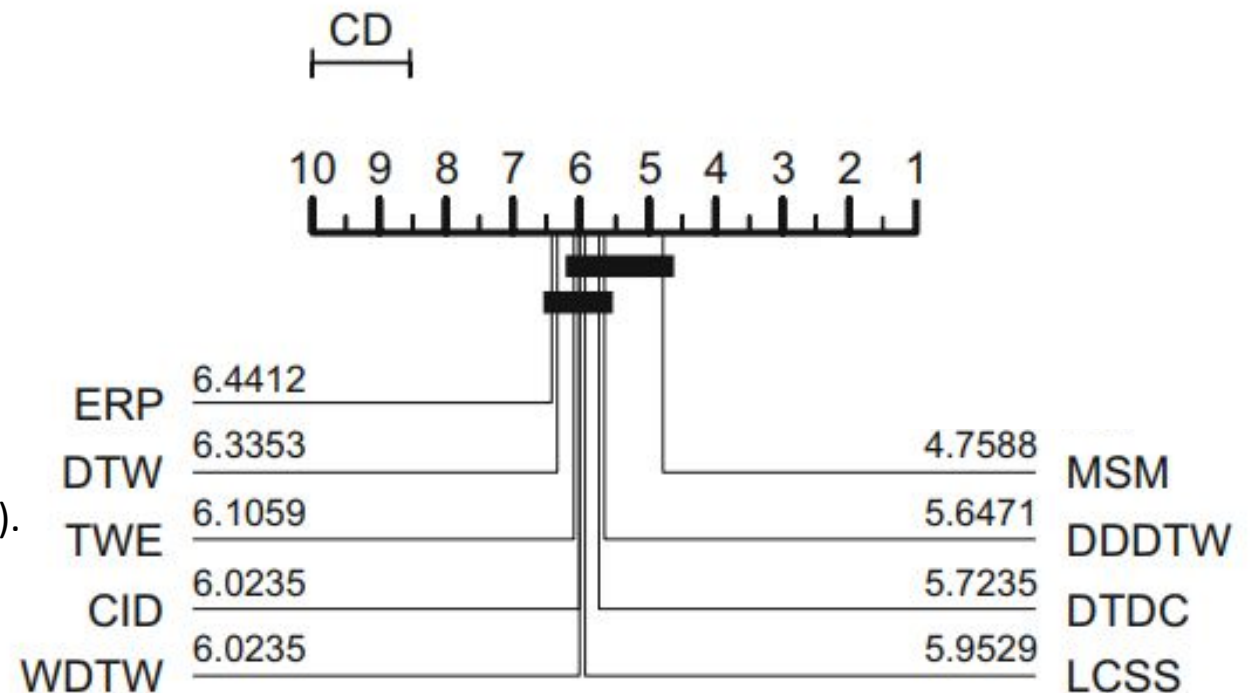# Time series classification



Corn

Wheat

Soybean

Broad-Leaved Tree

# Research into time series classification has accelerated very quickly over the last 5 years

Until recently,

- Many specialised time series classifiers developed

- But none dominated on accuracy on the UCR repository (85 datasets)
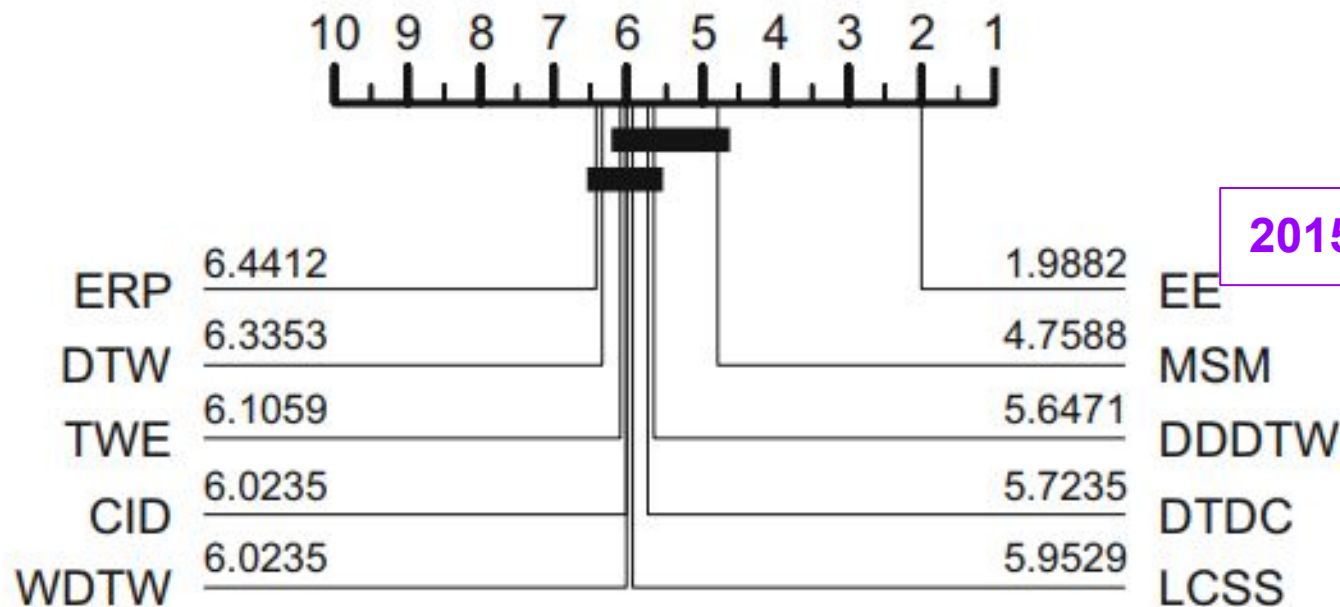
Bagnall, A., Lines, J., Bostrom, A., Large, J., & Keogh, E. (2017). The great time series classification bake off: a review and experimental evaluation of recent algorithmic advances. *Data Mining and Knowledge Discovery*, 31(3), 606-660.

| | | |
|---|---|---|
| ERP | 6.4412 | |
| DTW | 6.3353 | 4.7588 MSM |
| TWE | 6.1059 | 5.6471 DDDTW |
| CID | 6.0235 | 5.7235 DTDC |
| WDTW | 6.0235 | 5.9529 LCSS |

# A revolution in time series classification

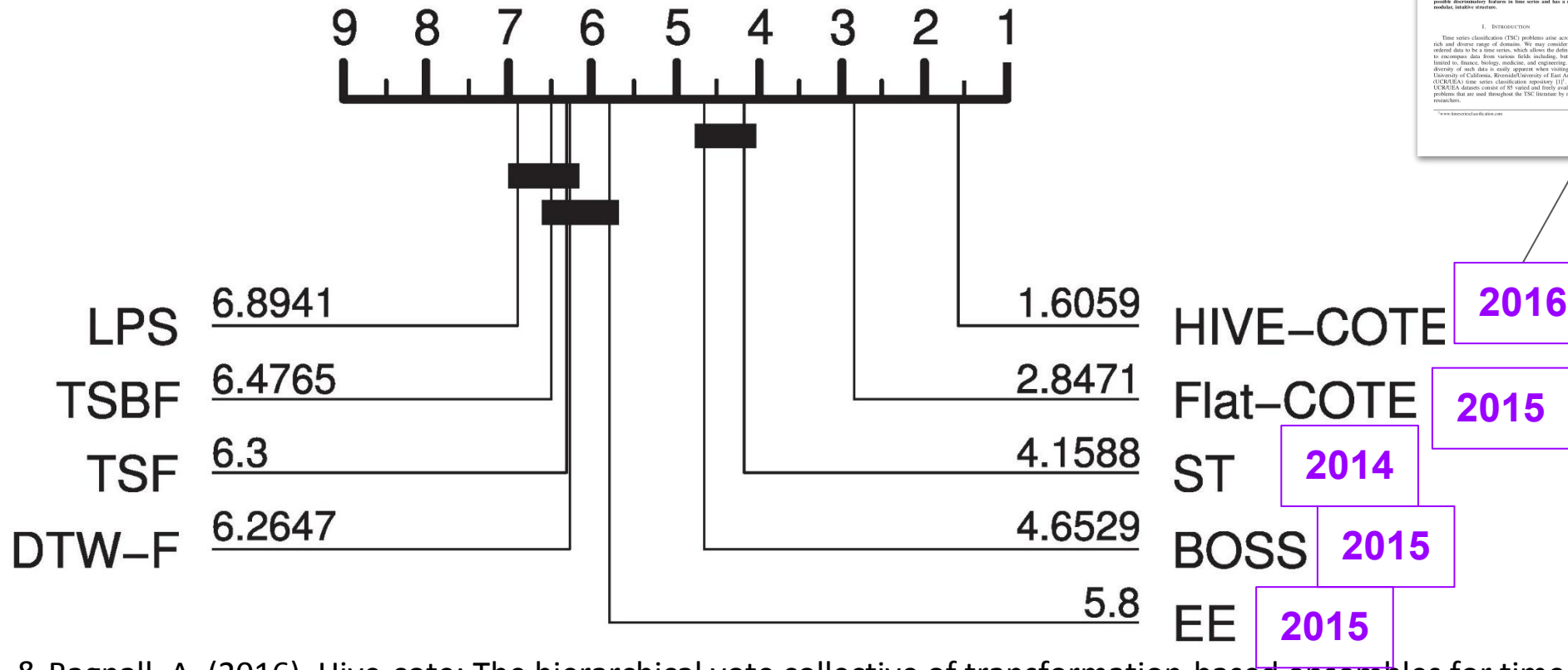Ensembles* have swept all before them!

* i.e. Tony, Jason, James and Aaron :)



Lines, J. & Bagnall, A., Time Series Classification with Ensembles of Elastic Distance Measures, *Data Mining and Knowledge Discovery*, 2015.
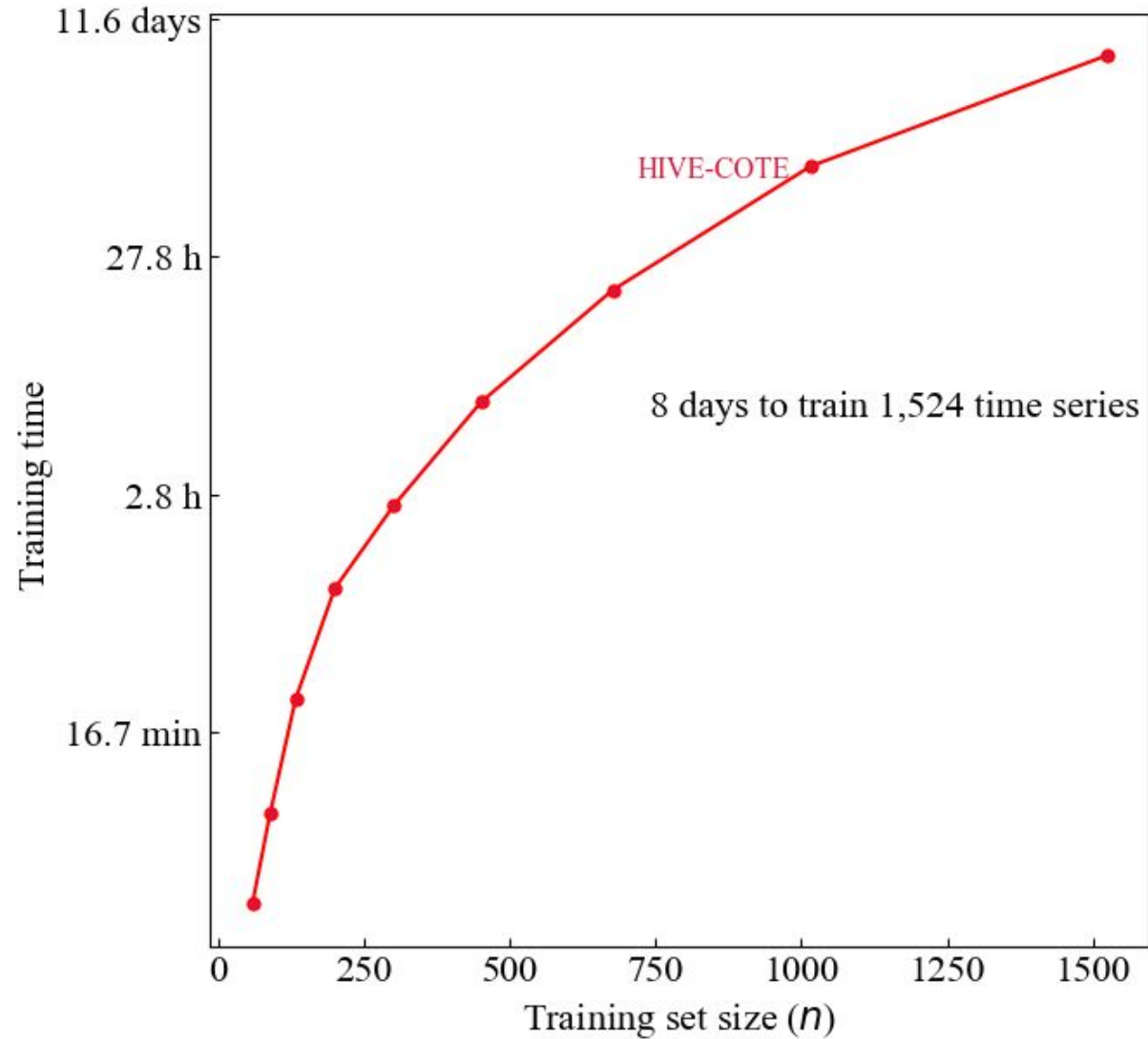
# A revolution in time series classification

They didn't stop there: a leap forward around 2015



Lines, J., Taylor, S., & Bagnall, A. (2016). Hive-cote: The hierarchical vote collective of transformation-based ensembles for time series classification. In *2016 IEEE 16th international conference on data mining (ICDM)* (pp. 1041-1046). IEEE.

# However, the most accurate ensembles do not scale

# State-of-the-art methods

# Training time complexity

Historical baseline

- ▶ 1-NN with Dynamic Time Warping (DTW)
  [Ratanamahatana and Keogh, 2005 ; Ding *et al.*, 2008]
- ▶ window size set by cross-validation

$$O(n^2 l^3)$$

Four leading classification algorithms  [Bagnall *et al.*, 2017]

- ▶ Bag-Of-SFA-Symbols (BOSS)  [Schäfer, 2015]
- ▶ Shapelet Transform (ST)  [Hills *et al.*, 2014]
- ▶ Elastic Ensembles (EE)  [Lines and Bagnall, 2015]
- ▶ Collective Of Transformation-based Ensembles (COTE)
  [Bagnall *et al.*, 2015]

$$O(n^2 l^3)$$
$$O(n^2 l^4)$$
$$O(n^2 l^3)$$

lower bounded by EE and ST algorithms

For 1M training instances, training the EE algorithm would require 73,000 days, **200 years**!

# However, the most accurate ensembles do not scale

A. Bagnall et al.

Overall, our results indicate that COTE is, on average, clearly superior to other published techniques. It is on average 8% more accurate than DTW. However, COTE is a starting point rather than a final solution. Firstly, the no free lunch theorem leads us to believe that no classifier will dominate all others. The research issues of most interest are what types of algorithm work best on what types of problem and can we tell *a priori* which algorithm will be best for a specific problem. Secondly, COTE is hugely computationally intensive. It is trivial to parallelise, but its run time complexity is bounded by the Shapelet Transform, which is $O(n^2m^4)$ and the parameter searches for the elastic distance measures, some of which are $O(n^3)$. ST and EE are also trivial to distribute, but there is a limit to the number of processors anyone can run in parallel. An algorithm that is faster than COTE but not significantly less accurate would be a genuine advance in the field. Finally, we are only looking at a very restricted type of

Bagnall, A., Lines, J., Bostrom, A., Large, J., & Keogh, E. (2017). The great time series classification bake off: a review and experimental evaluation of recent algorithmic advances. *Data Mining and Knowledge Discovery*, *31*(3), 606-660.
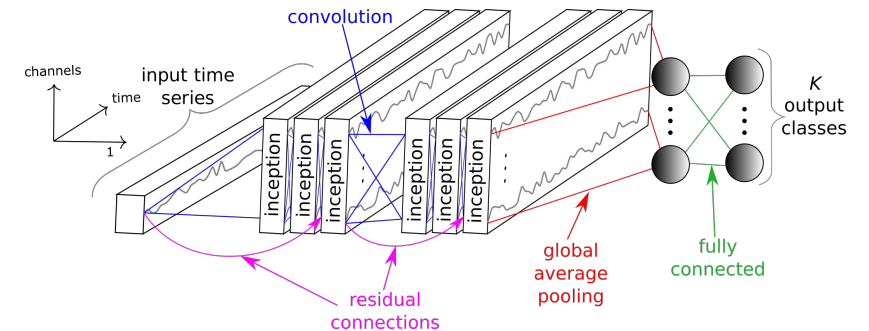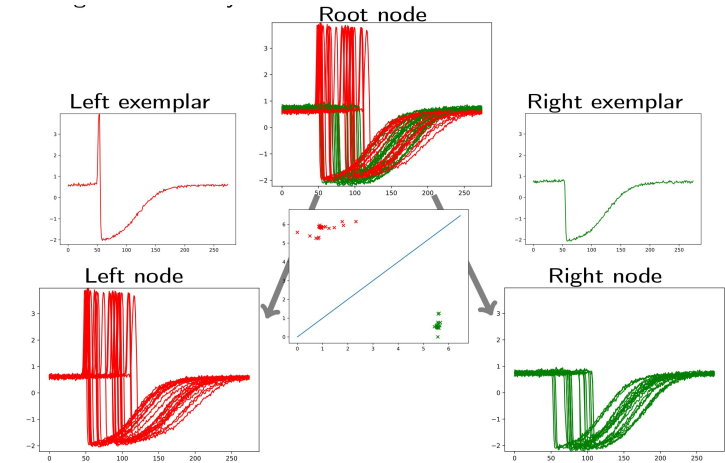
# Talk Outline



Highly accurate and scalable TSCs

- Tree-based: *Proximity Forest* and *TS-CHIEF*
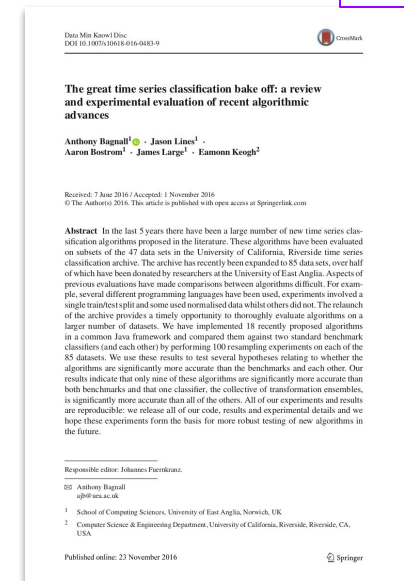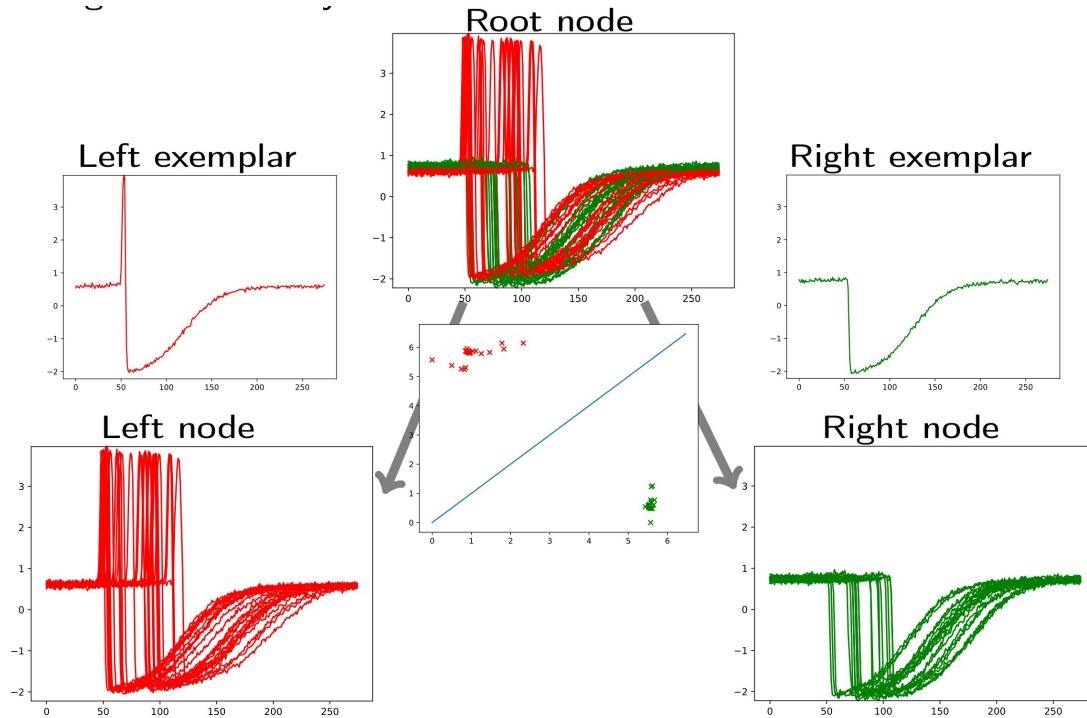
- Deep Learning: *InceptionTime*

This talk is super fresh!

- 1 DAMI 2019 paper

- 2 arxiv papers submitted in the last 3 months
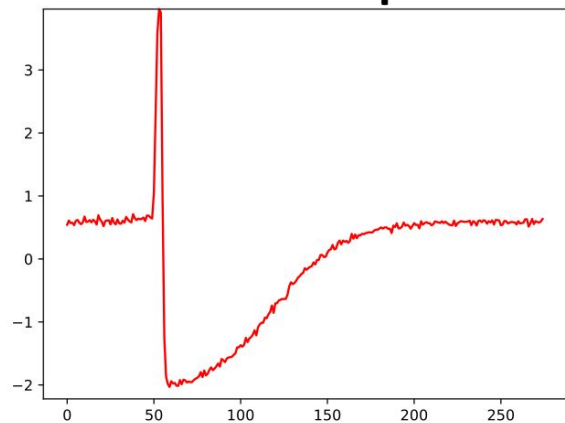
# Part 1: Proximity Forest (PF)



B. Lucas, A. Shifaz, C. Pelletier, L. O'Neill, N. Zaidi, B. Goethals, F. Petitjean, G. Webb (2019). Proximity Forest: An effective and scalable distance-based classifier for time series. *Data Mining and Knowledge Discovery*, *33*(3), 607-635.
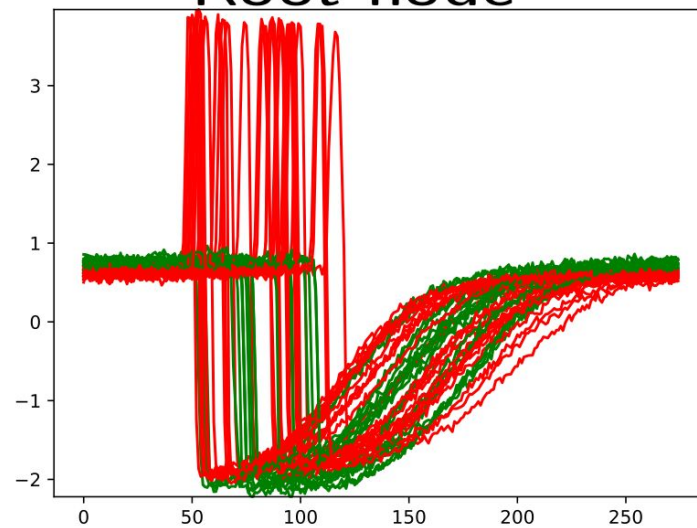
http://bit.ly/ProximityForest

# Proximity Forest

Starting point: How to make Elastic Ensemble (EE) scalable?

- We need a divide-and-conquer approach to be efficient

- We want to emulate Elastic Ensemble as closely as possible to allow clear comparison of fundamental strategies

- But tree-base splits don't work for time series because no attribute/value representation
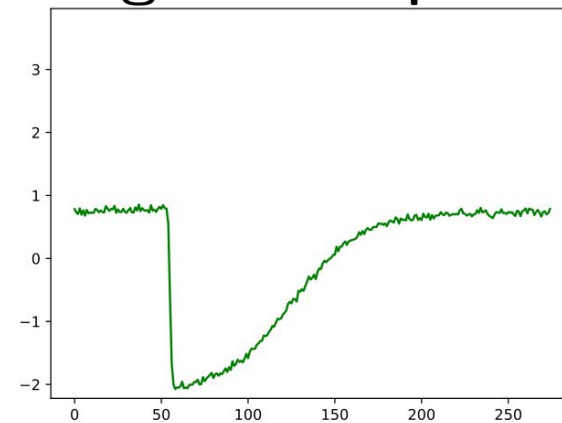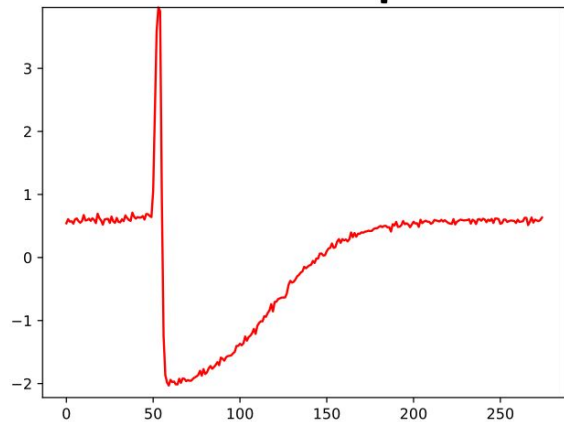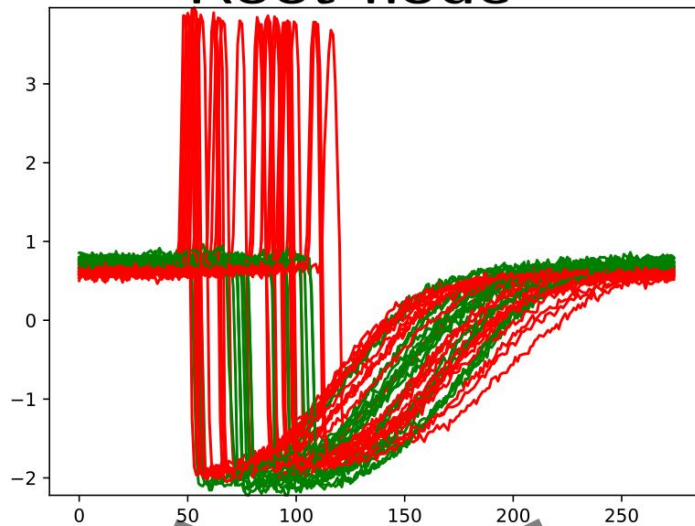
**Left exemplar**

**Root node**

**Right exemplar**

Root node

Left exemplar
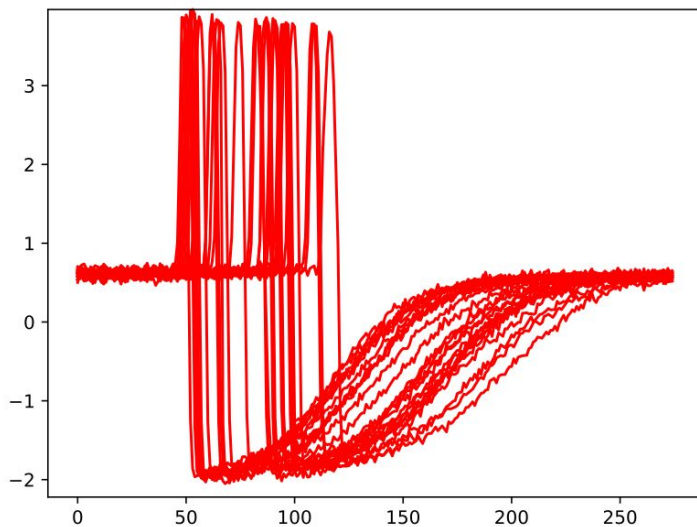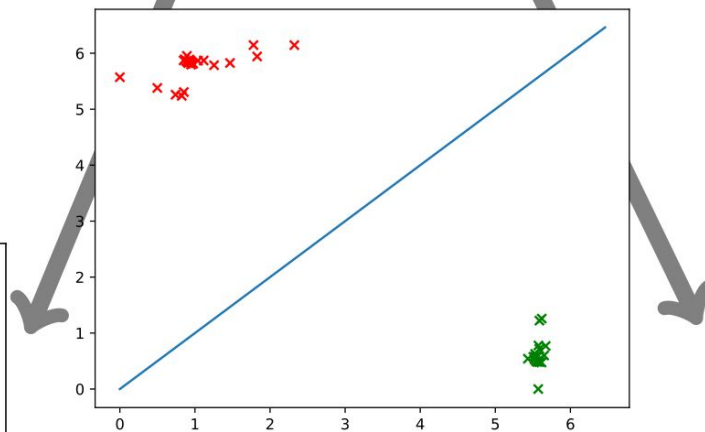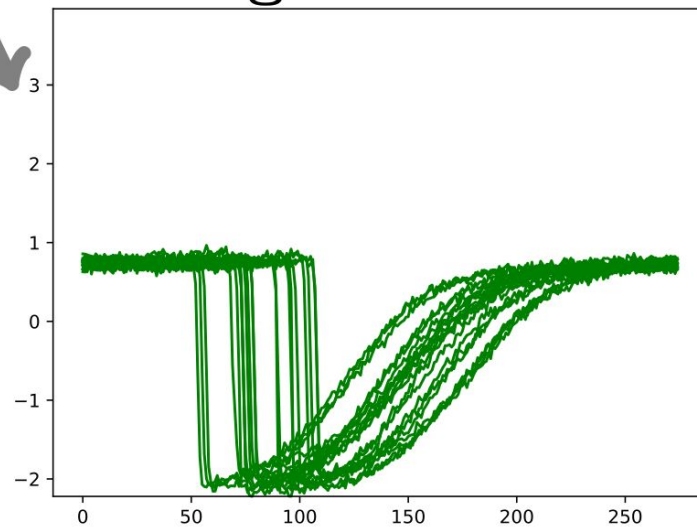
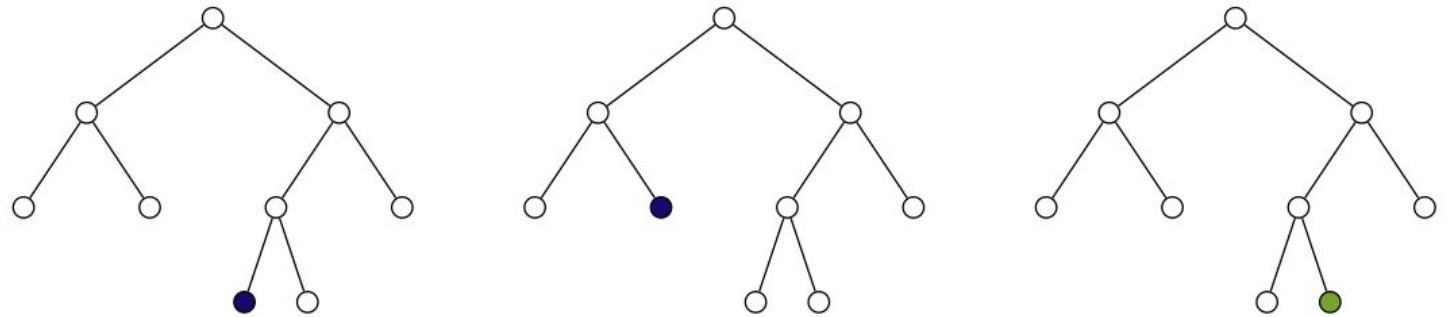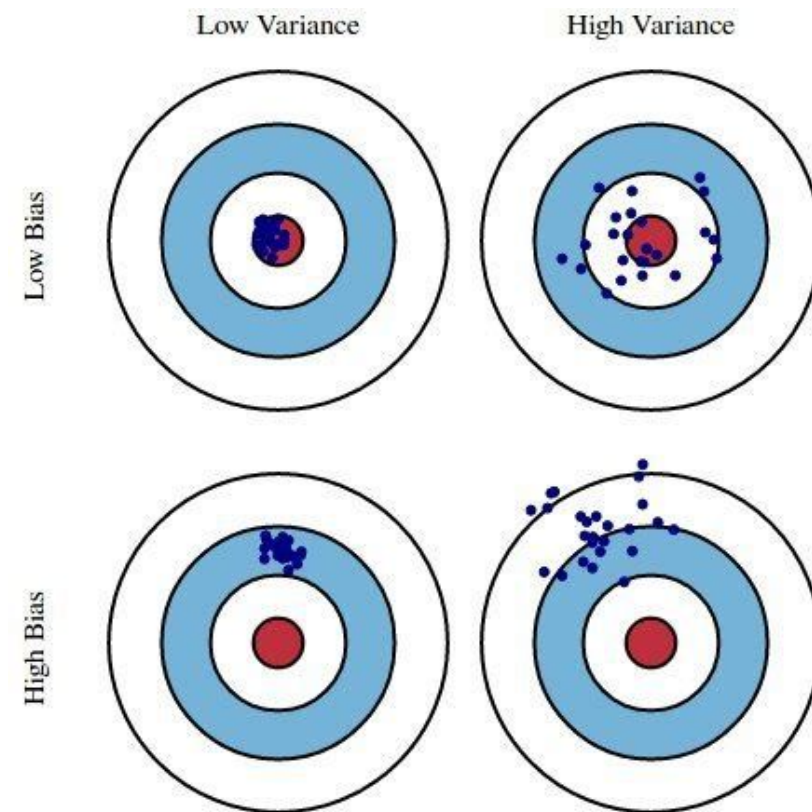Right exemplar

Left node

Right node

# Proximity Tree

- Replace conventional decision tree splits with similarity comparisons using specialised time series methods
  → Makes the most of 40 years of research into designing appropriate measures for time series (DTW, TWE, MSM, LCSS, etc)

- Each branch has an exemplar associated with it

- One exemplar per class

- Each split in the tree has (1) a measure and (2) a parametrization

- For classification, series $S$ p
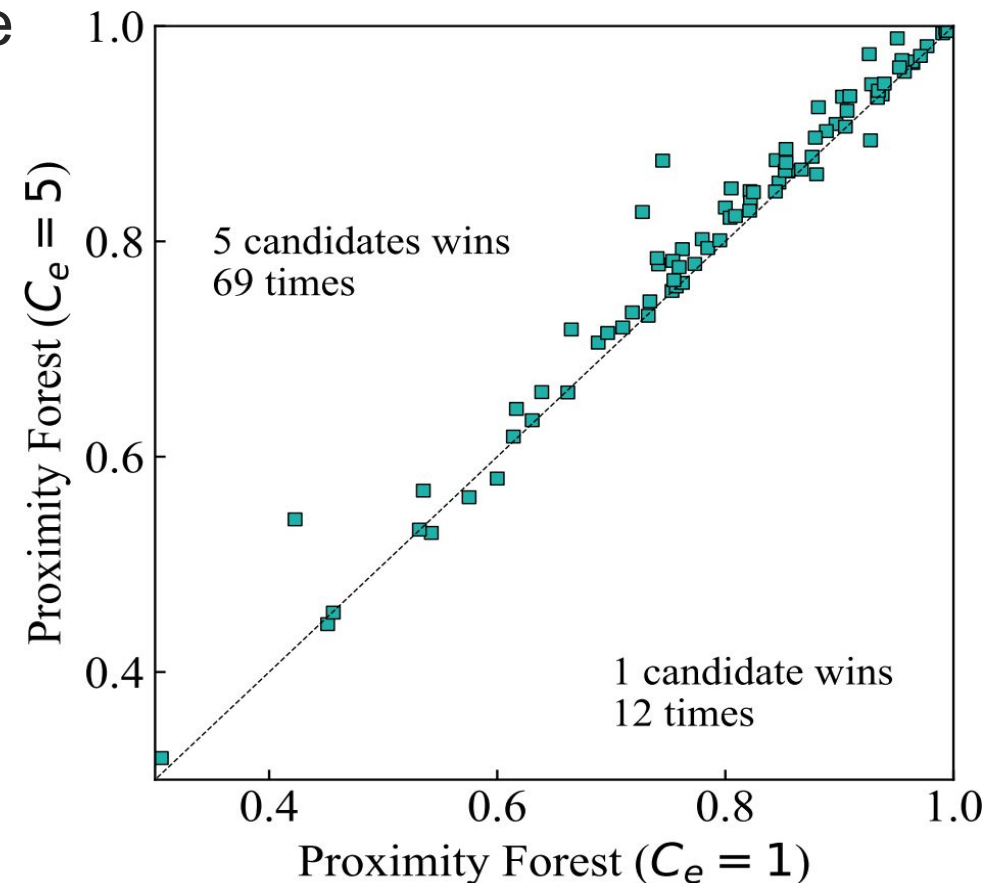  most similar

# Stochastic choices for speed and diversity

- Exemplars chosen at random among series at the node

- Distance measures and their parameterizations chosen at random from those used by EE

- Random choices have low bias and ensembling removes the resulting variance

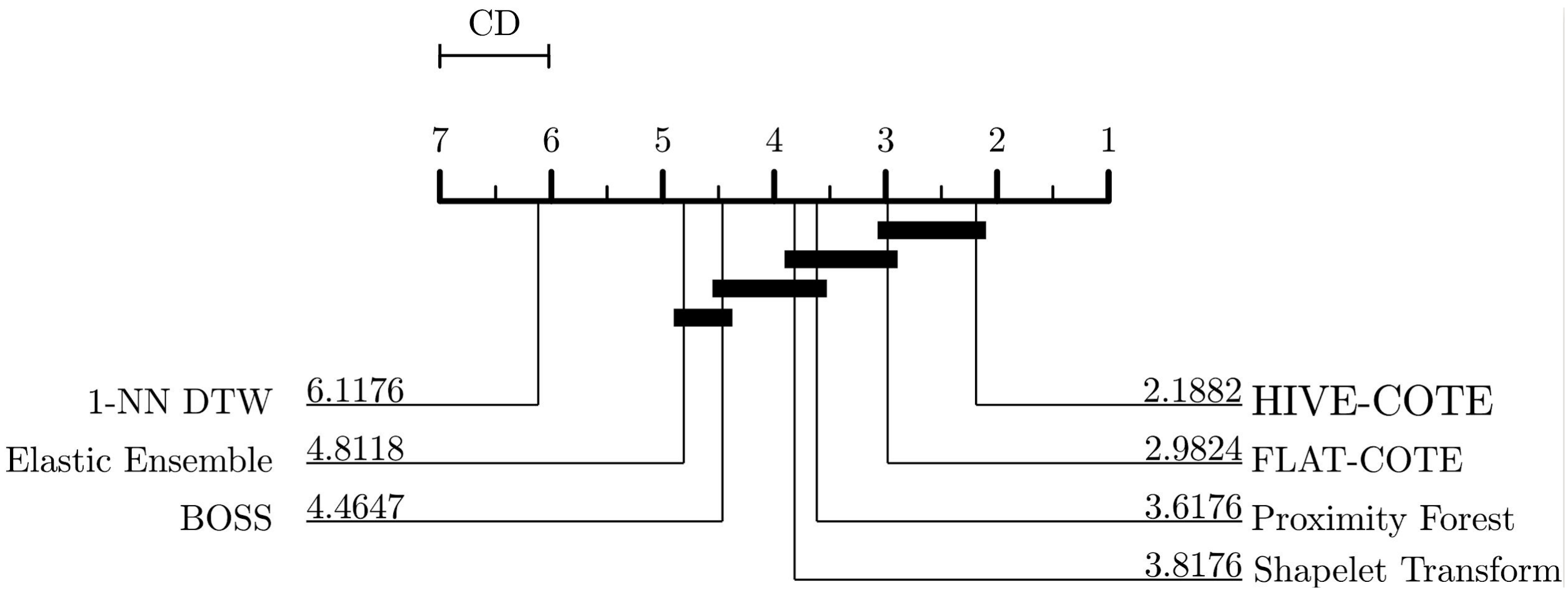- The major training time cost is passing training examples down the tree



$$\text{MSE}(H) = \overline{bias}(H)^2 + \frac{1}{|H|}\overline{variance}(H) + \left(1 - \frac{1}{|H|}\right)\overline{covariance}(H)$$

N. Ueda and R. Nakano, "Generalization error of ensemble estimators," *Proceedings of International Conference on Neural Networks (ICNN'96)*, Washington, DC, USA, 1996, pp. 90-95 vol.1. doi: 10.1109/ICNN.1996.548872
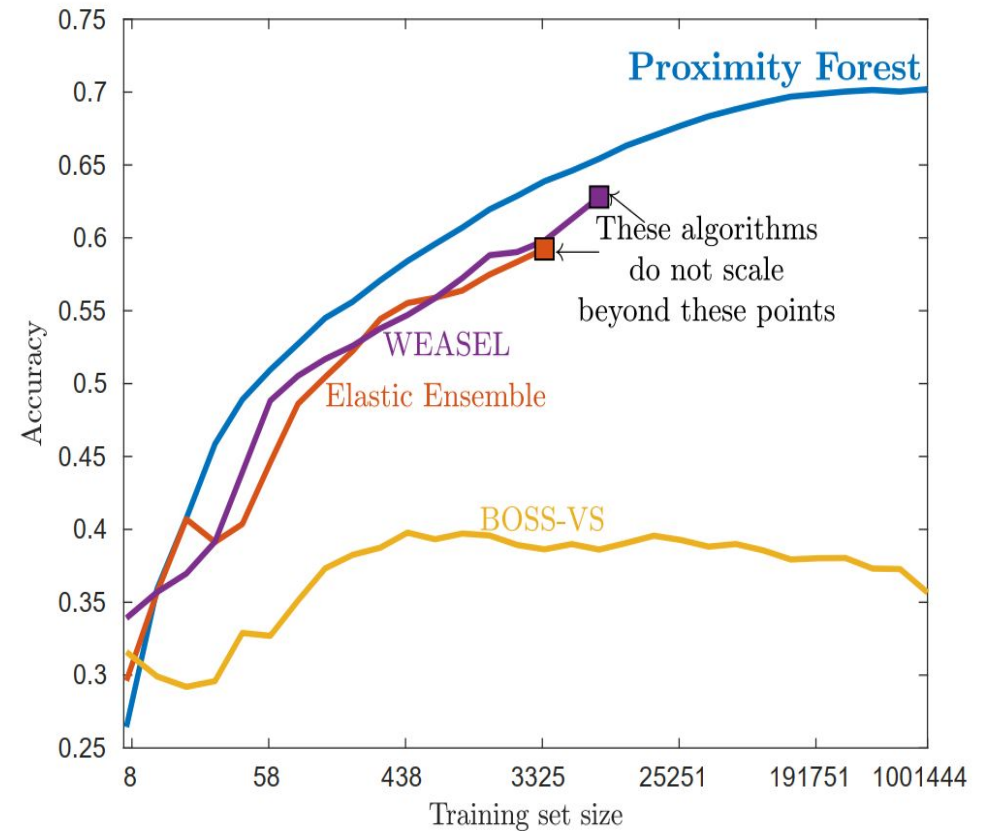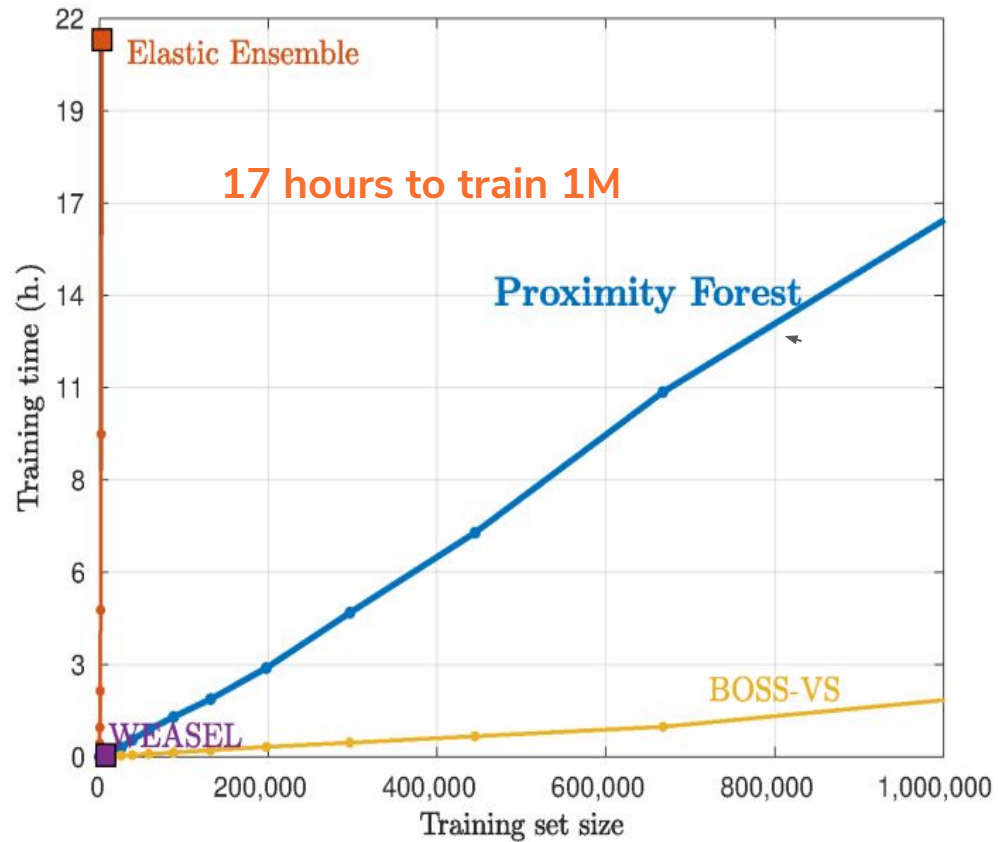
# Select between multiple random candidates at each node

- Use GINI to select best from five candidate splits

- Increases covariance, decreases variance

  → so we don't need too many trees

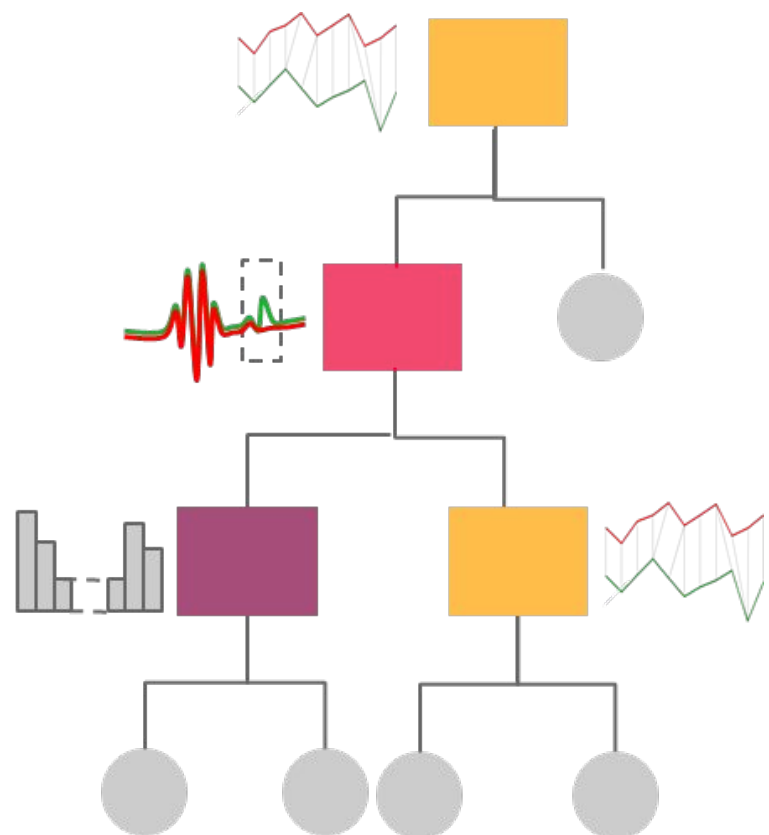  → faster training

  → faster classification



5 candidates wins 69 times

1 candidate wins 12 times

Proximity Forest ($C_e = 5$) vs Proximity Forest ($C_e = 1$)

# Scalability evaluated on 1M instances of Satellite Image Time Series (SITS) dataset

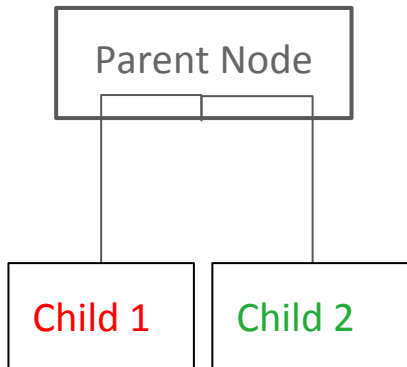# Part 2: TS-CHIEF



2019 - arxiv

http://bit.ly/TS-CHIEF

A. Shifaz, C. Pelletier, F. Petitjean and G. Webb (2019). TS-CHIEF: A Scalable and Accurate Forest Algorithm for Time Series Classification. *under review.* *https://arxiv.org/abs/1906.10329*
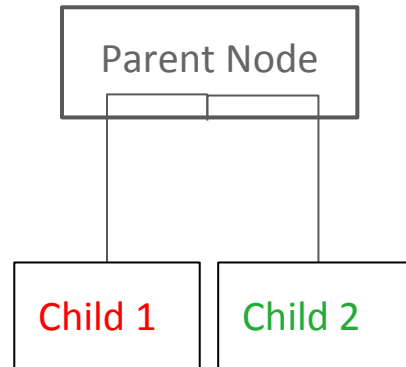
# Time Series Combination of Integrated Embeddings Forest (TS-CHIEF)

Similarity-based
(Proximity Forest)

Dictionary-based

Interval-based
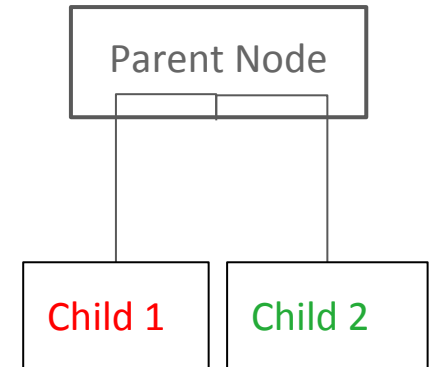


Candidate split 1

Candidate split 2

Candidate split 3

- Candidates selected at random from all three strategies
- Selection using Gini Index

# Time Series Combination of Heterogeneous Integrated Embeddings Forest (TS-CHIEF)
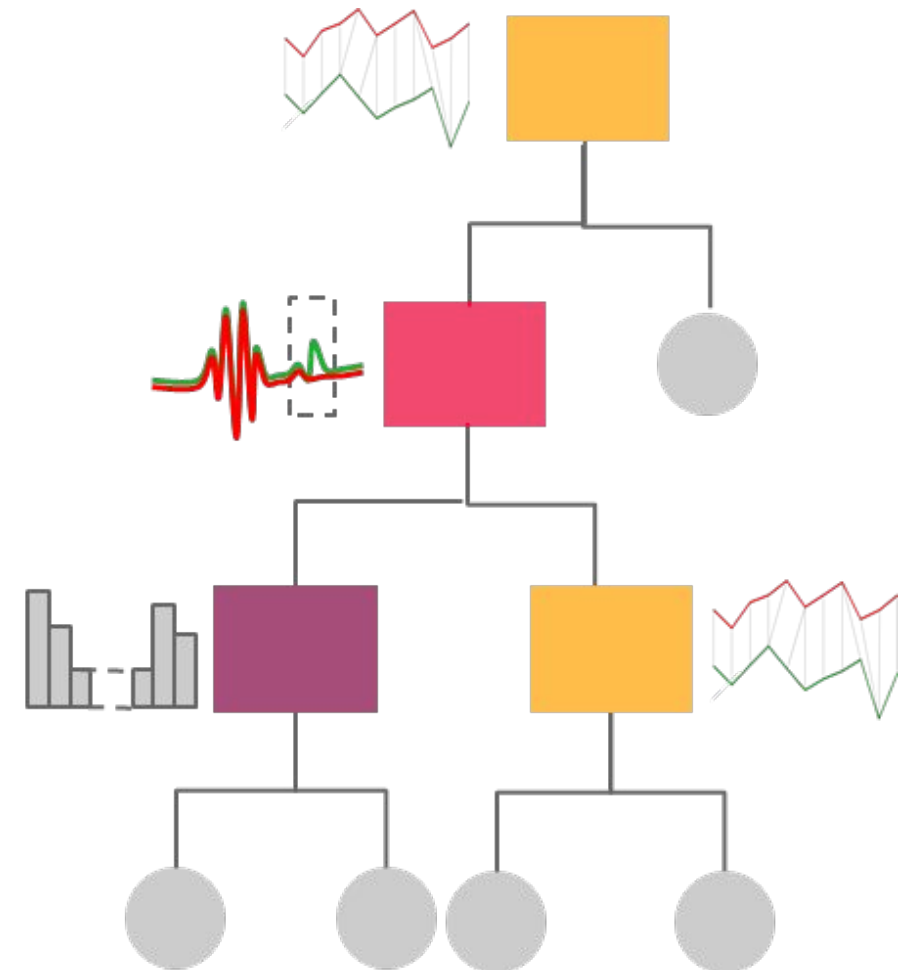
- TS−CHIEF trees combine three splitting functions
- Candidate splits selected at random
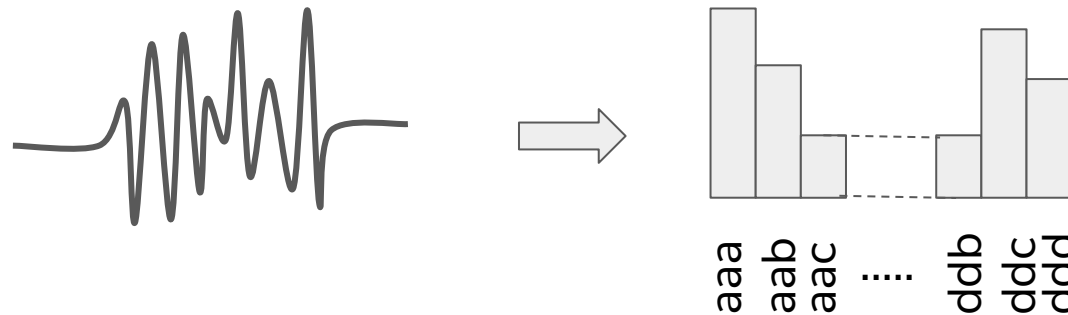- Final selection using Gini Index

# TS-CHIEF : Dictionary-based splitter

- Precomputes a pool of BOSS transformations at forest level



- At node select a random transformation

- At node selects reference histograms per class (exemplars)

- Uses histogram similarity measure

- Partitions the data based on the proximity to reference histograms

- Original BOSS: Uses cross validation

- TS-CHIEF: Uses random transformations

P. Schäfer (2015), The BOSS is concerned with time series classification in the presence of noise, *Data Mining and Knowledge Discovery*, Vol 29, Num 6, pp 1505–1530.
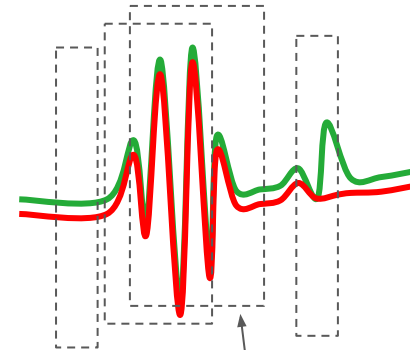
# TS-CHIEF : Interval-based splitter

- Select random intervals and transforms
  - time (ACF, PACF, AR) and frequency (PS)

- Attribute-value split similar to classic decision tree

At tree level: RISE
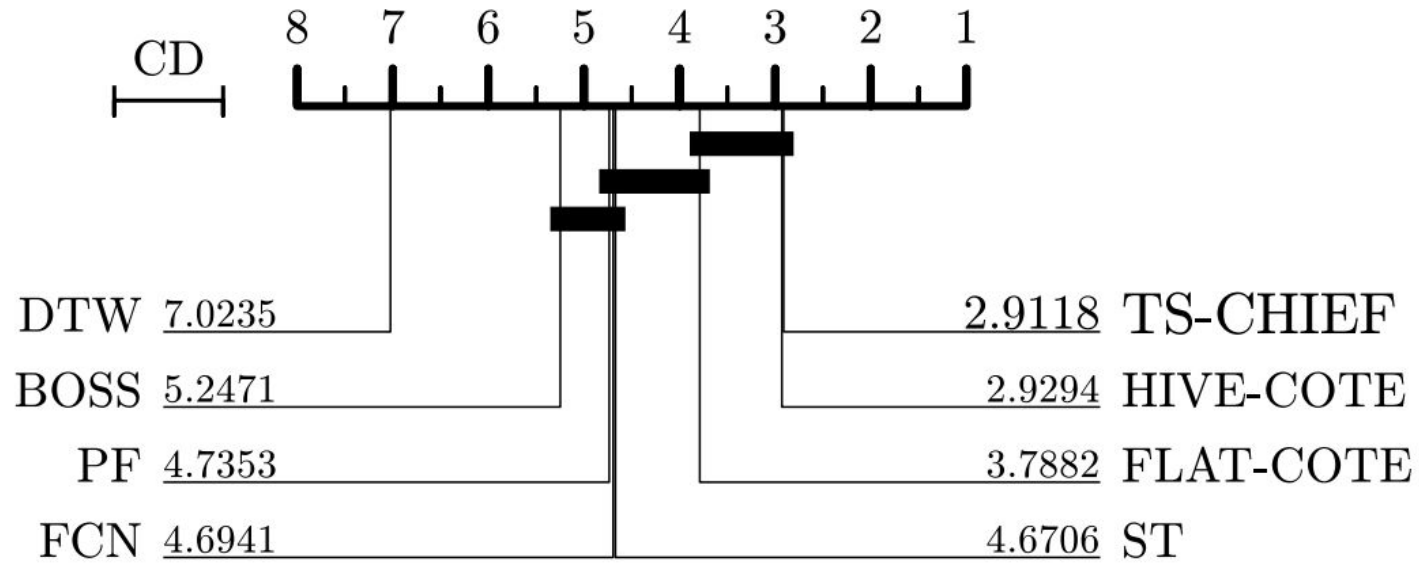
At node level: TS-CHIEF



Intervals selected on global discrimination ability

random intervals

# Accuracy on 85 UCR datasets

# Accuracy on 85 UCR datasets



**TS-CHIEF vs HIVE-COTE**

41 wins
35 loss
9 ties
p=0.42 (Wilcoxon's test)

# Accuracy on 85 UCR datasets



| | | | |
|---|---|---|---|
| DTW 7.0235 | | | 2.9118 TS-CHIEF |
| BOSS 5.2471 | | | 2.9294 HIVE-COTE |
| PF 4.7353 | | | 3.7882 FLAT-COTE |
| FCN 4.6941 | | | 4.6706 ST |



**TS-CHIEF vs HIVE-COTE**

41 wins
35 loss
9 ties
p=0.42 (Wilcoxon's test)

Training time vs training size

| Training Size | HIVE-COTE | TS-CHIEF |
|---|---|---|
| 1,500 time series | 8 days | 13 min (900x faster) |
| 130,000 time series | 230 years (estimated) | 2 days (46,000x faster) |

# Part 3: InceptionTime

H. Ismail Fawaz, B. Lucas, G. Forestier, C. Pelletier, D. Schmidt, J. Weber, G. Webb, L. Idoumghar, P-A. Muller, F. Petitjean (2019). InceptionTime: Finding AlexNet for Time Series Classification. *under review.* https://arxiv.org/abs/1909.04939

# Deep Learning

- Revolutionized the field of computer vision [1]

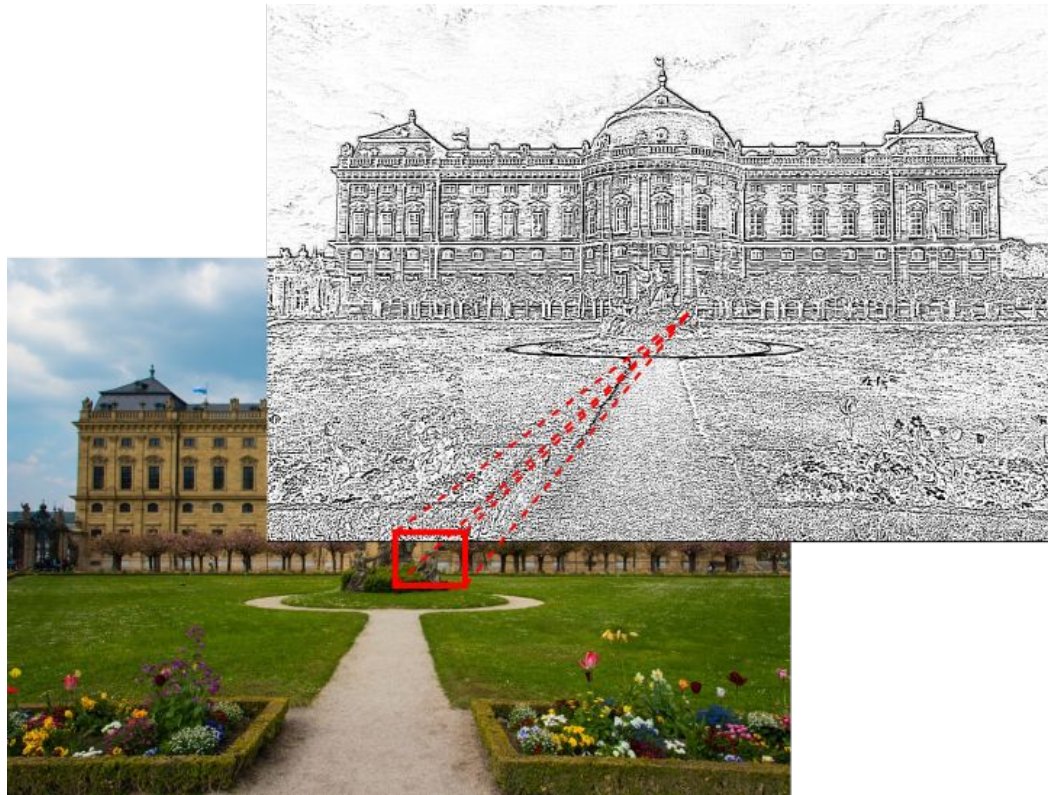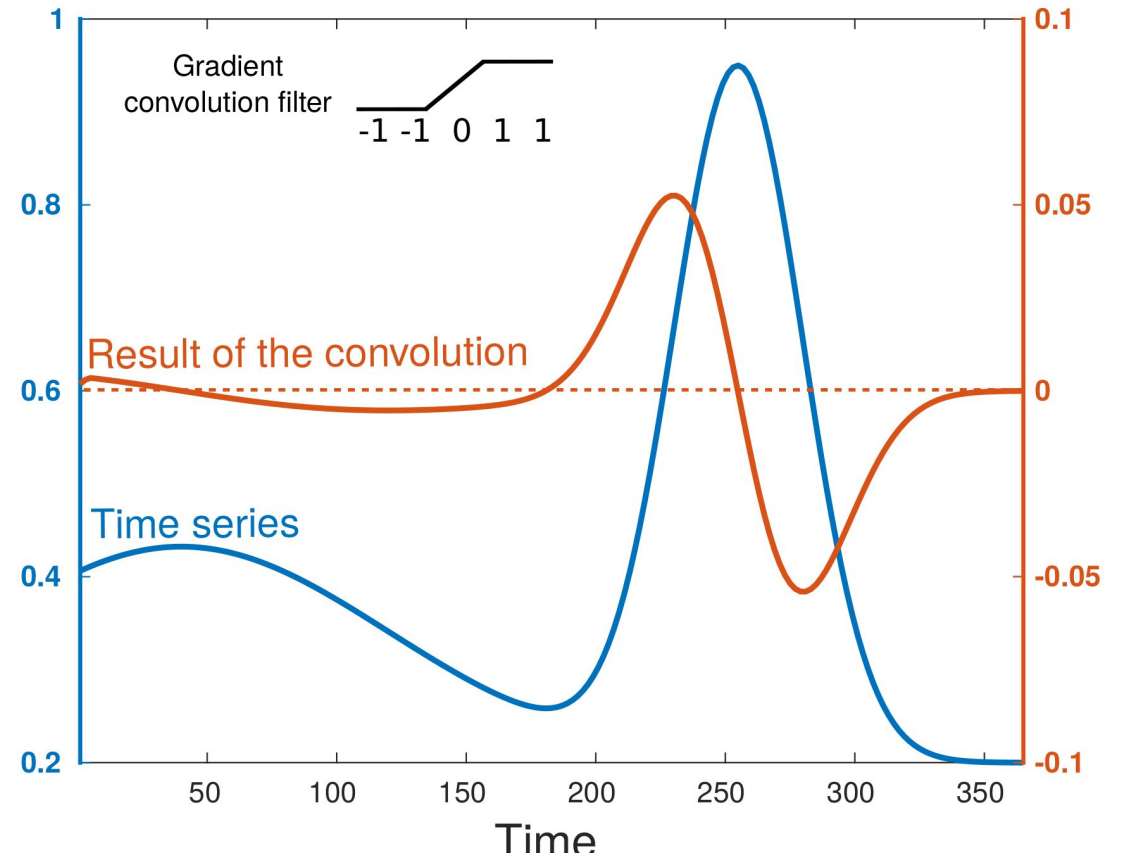- Reached human level performance in image recognition tasks [2]

- Adopted by the Natural Language Processing (NLP) community [3]

- Improved state of the art speech recognition systems[4]

1. Krizhevsky, A., Sutskever, I., & Hinton, G. E. (2012). Imagenet classification with deep convolutional neural networks. In Advances in Neural Information Processing Systems
2. Szegedy, C., Liu, W., Jia, Y., Sermanet, P., Reed, S., Anguelov, D., et al. (2015). Going deeper with convolutions. In Proceedings of IEEE Conference on Computer Vision and Pattern Recognition
3. Young, T., Hazarika, D., Poria, S., & Cambria, E. (2018). Recent trends in deep learning based natural language processing. IEEE Computational intelligenCe magazine, 13(3), 55-75.
4. Hinton, G., Deng, L., Yu, D., Dahl, G. E., Mohamed, A. R., Jaitly, N., et al. (2012). Deep neural networks for acoustic modeling in speech recognition: The shared views of four research groups. IEEE Signal Processing Magazine

# Convolution on images vs. time series



The result of a applying an edge detection convolution on an image



Gradient convolution filter
-1 -1 0 1 1

Result of the convolution

Time series

Time

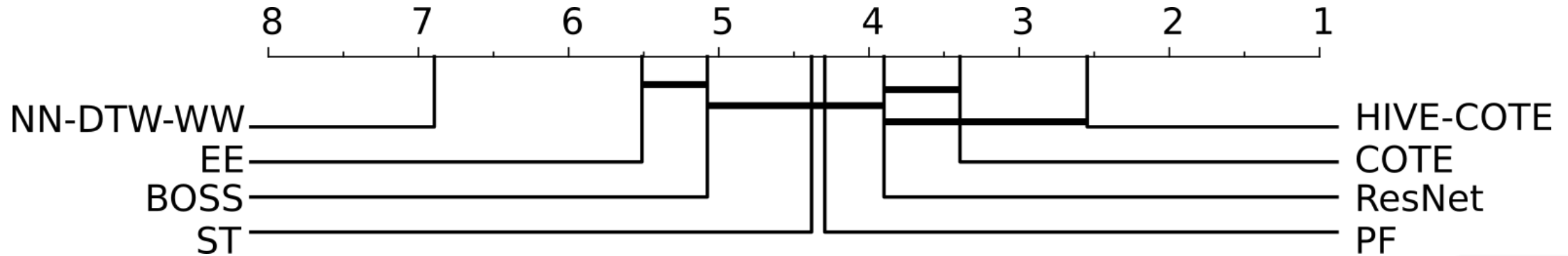# Convolution on images vs. time series



The result of a applying an edge detection convolution on an image



The result of applying a learned discriminative convolution on the GunPoint dataset

# Deep learning for Time Series Classification



A critical difference diagram showing how ResNet still lacks behind the state of the art classifiers [1]

- Residual Network (ResNet) was originally proposed in [2]
- Currently is the state-of-the-art deep learning model for TSC [1]
- Designed to be a "baseline architecture" for TSC

1. Ismail Fawaz, H., Forestier, G., Weber, J., Idoumghar, L., & Muller, P. A. (2019). Deep learning for time series classification: a review. Data Mining and Knowledge Discovery, 33(4), 917-963.
2. Wang, Z., Yan, W., & Oates, T. (2017, May). Time series classification from scratch with deep neural networks: A strong baseline. In IEEE *International Joint Conference on Neural Networks.*
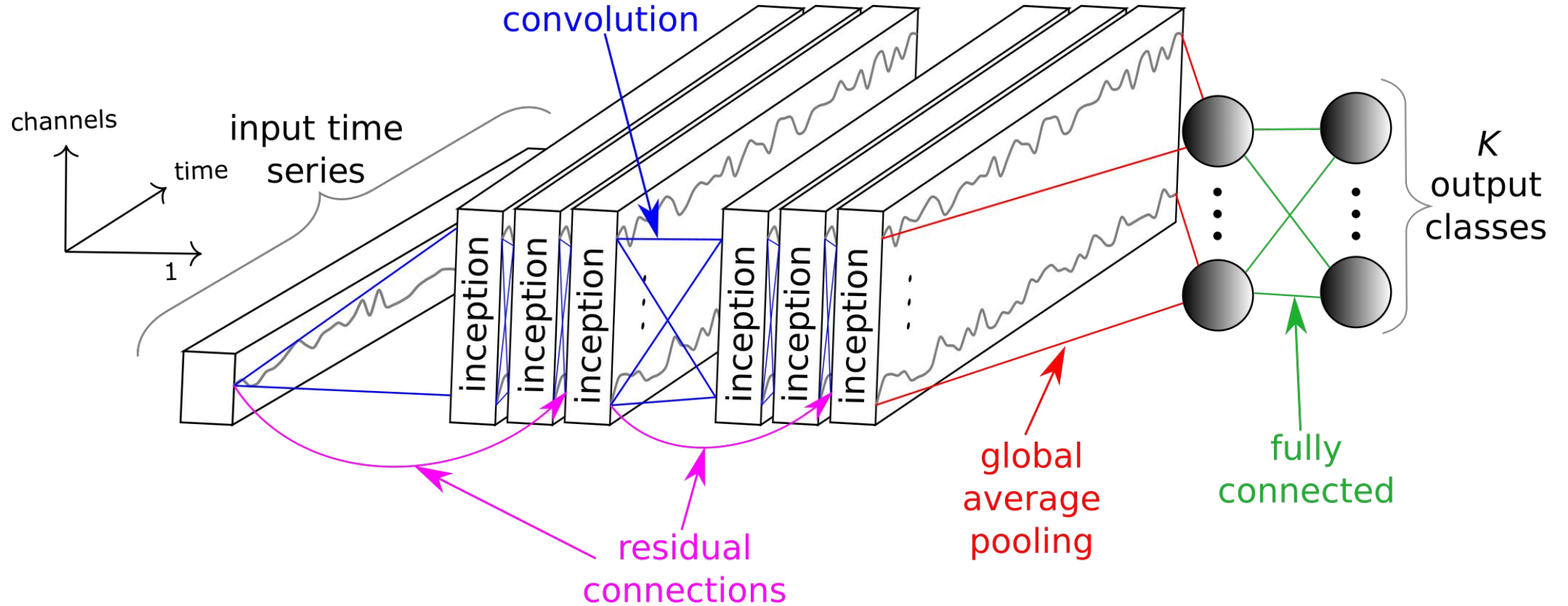
# Inception

- Originally proposed by Google for image recognition problems [1]

- Further developed to reach state-of-the-art results on ImageNet [2]

- Main idea:
  - Apply convolutions of different resolutions to capture different patterns
  - Use a bottleneck layer in order to reduce the number of parameters

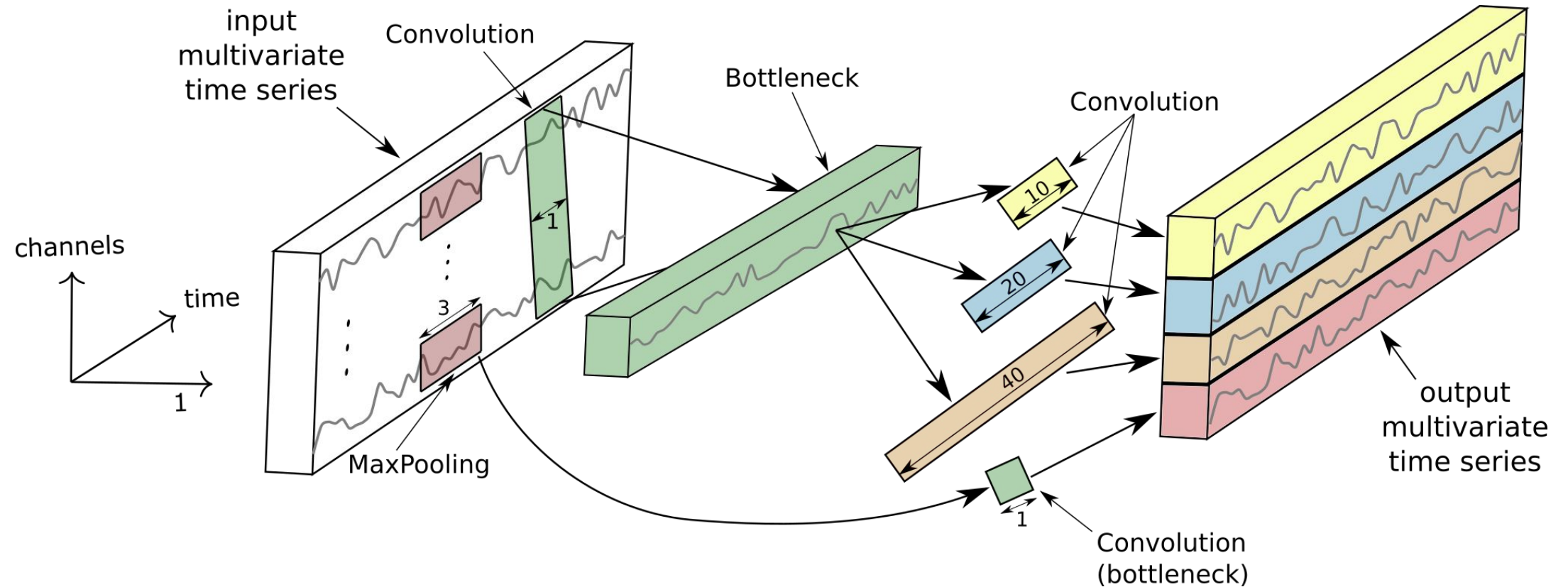- For TSC, Inception had not been yet explored

1. Szegedy, C., Liu, W., Jia, Y., Sermanet, P., Reed, S., Anguelov, D., Erhan, D., Vanhoucke, V., & Rabinovich, A. (2015). Going deeper with convolutions. In Proceedings of the IEEE conference on computer vision and pattern recognition (pp. 1-9).
2. Szegedy, C., Vanhoucke, V., Ioffe, S., Shlens, J., & Wojna, Z. (2016). Rethinking the inception architecture for computer vision. In *Proceedings of the IEEE conference on computer vision and pattern recognition* (pp. 2818-2826).
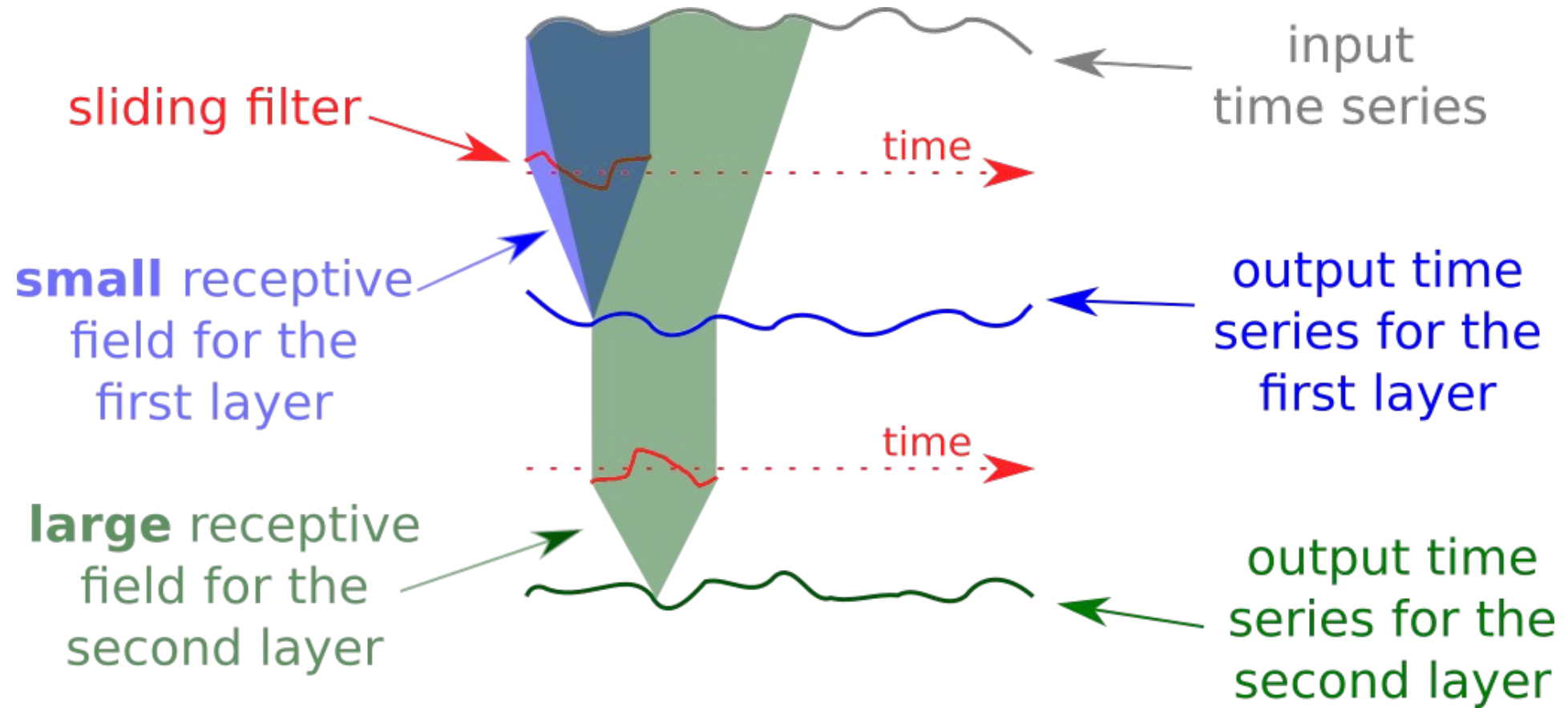
# Our InceptionTime architecture for TSC



Inception network for time series classification

# Inception **module** for time series classification



Inside our Inception module for time series classification

# Receptive Field (RF) of a neural network



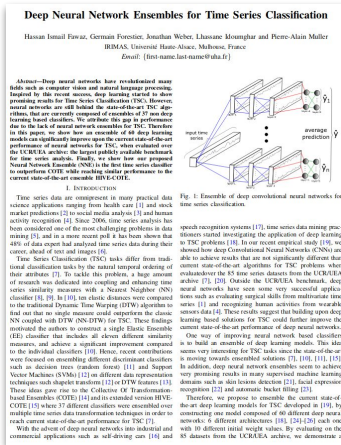Receptive field illustration for a two layers CNN

# InceptionTime: an ensemble of 5 networks

$$\hat{y}_{i,c} = \frac{1}{n} \sum_{j=1}^{n} \sigma_c(x_i, \theta_j) \quad | \quad \forall c \in [1, C]$$
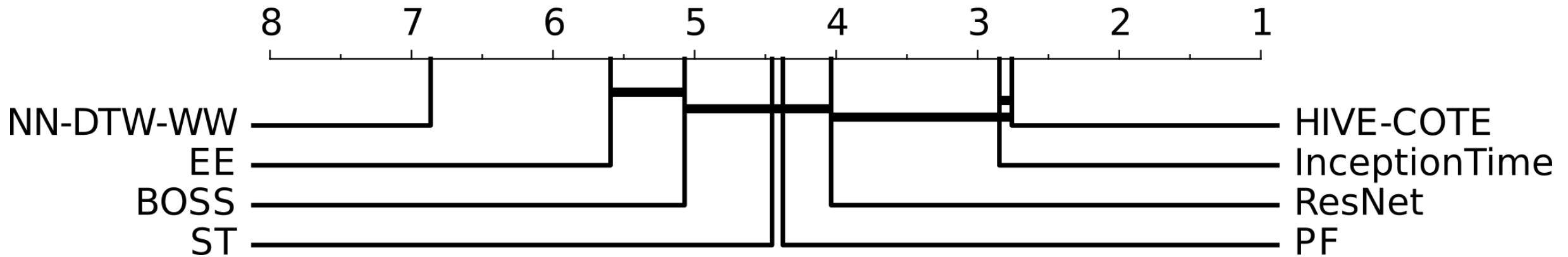
- Ensembling deep nets for TSC studied by Hassan in [1]
- Bias/variance tells us that this works because different initializations lead to very different networks (low covariance)

$$\text{MSE}(H) = \overline{bias}(H)^2 + \frac{1}{|H|} \overline{variance}(H) + \left(1 - \frac{1}{|H|}\right) \overline{covariance}(H)$$

1. Ismail Fawaz, H., Forestier, G., Weber, J., Idoumghar, L., & Muller, P. (2019). Deep neural network ensembles for time series classification. IEEE International Joint Conference on Neural Networks.

# Accuracy results on the UCR archive



Critical difference diagram showing the performance of InceptionTime compared to the current state-of-the-art classifiers of time series data
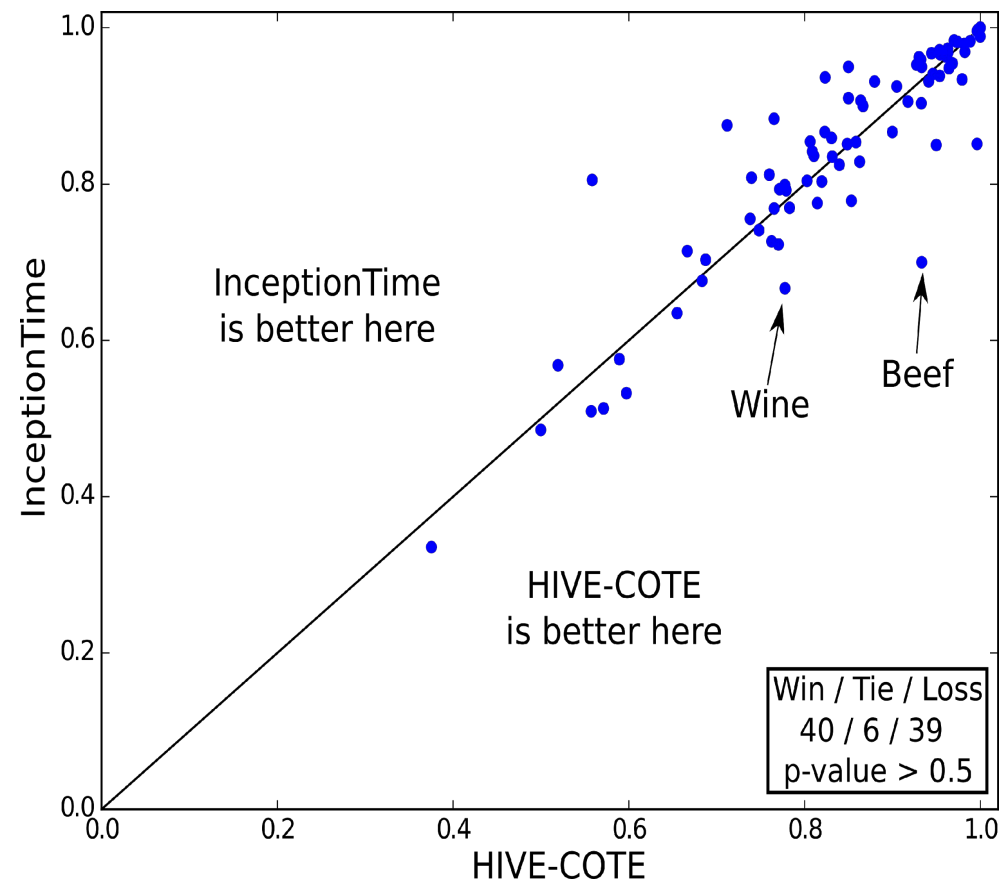
- InceptionTime reaches very similar results to HIVE-COTE

Bagnall, A., Lines, J., Bostrom, A., Large, J., & Keogh, E. (2017). The great time series classification bake off: a review and experimental evaluation of recent algorithmic advances. *Data Mining and Knowledge Discovery*, *31*(3), 606-660.
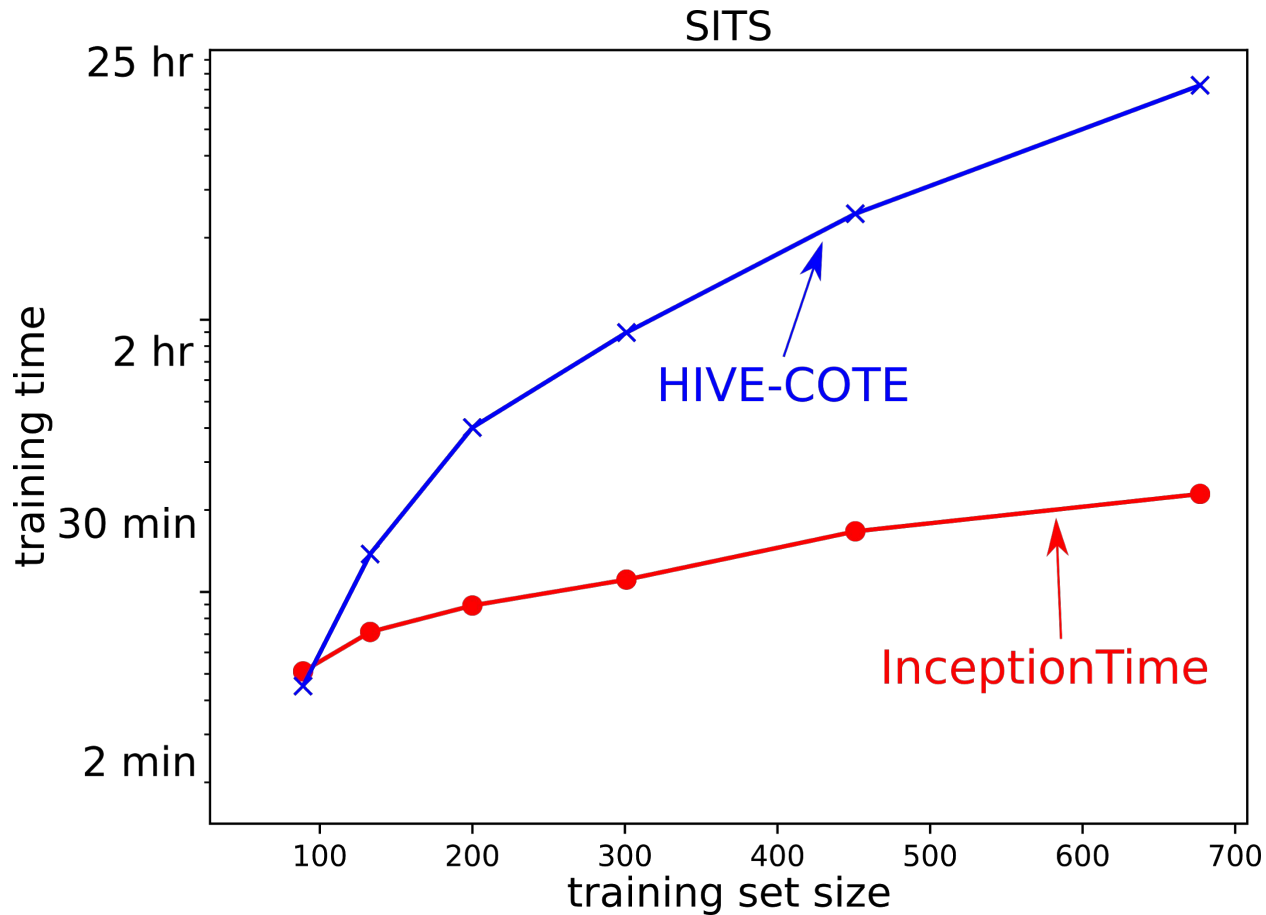
# Accuracy plot: InceptionTime vs HIVE-COTE

- InceptionTime is slightly better than HIVE-COTE on average [1]

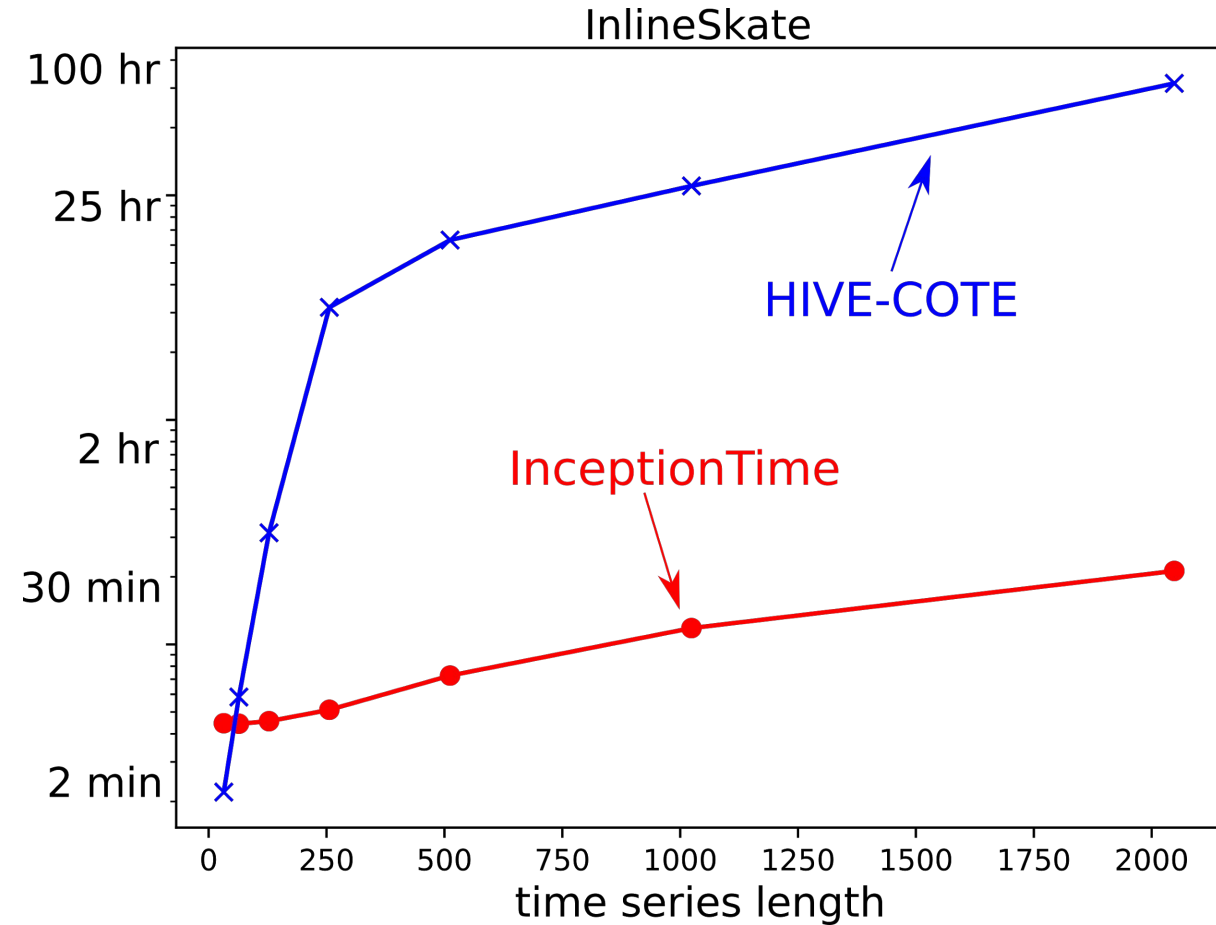- Wine and Beef were shown to benefit from transfer learning [2]

1. Lines, J., Taylor, S., & Bagnall, A. (2018). Time series classification with HIVE-COTE: The hierarchical vote collective of transformation-based ensembles. *ACM Transactions on Knowledge Discovery from Data (TKDD)*, *12*(5), 52.

2. Ismail Fawaz, H. I., Forestier, G., Weber, J., Idoumghar, L., & Muller, P. A. (2018). Transfer learning for time series classification. In IEEE International Conference on Big Data.
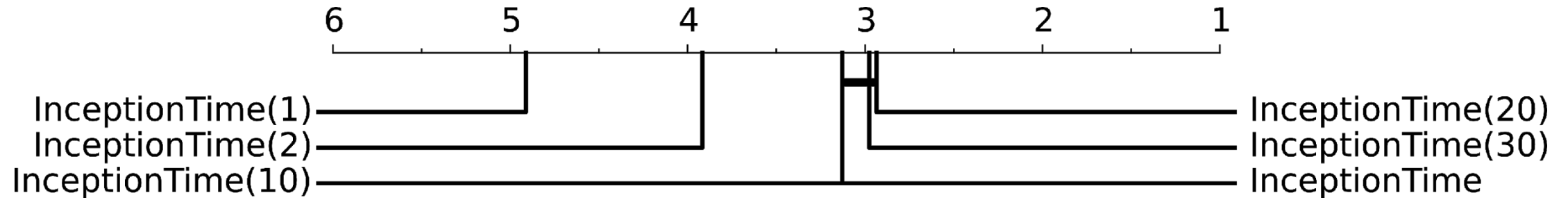
# Training time comparison with HIVE-COTE



Training time as a function of the training set size for the SITS dataset

Training time as a function of the series length for the InlineSkate dataset
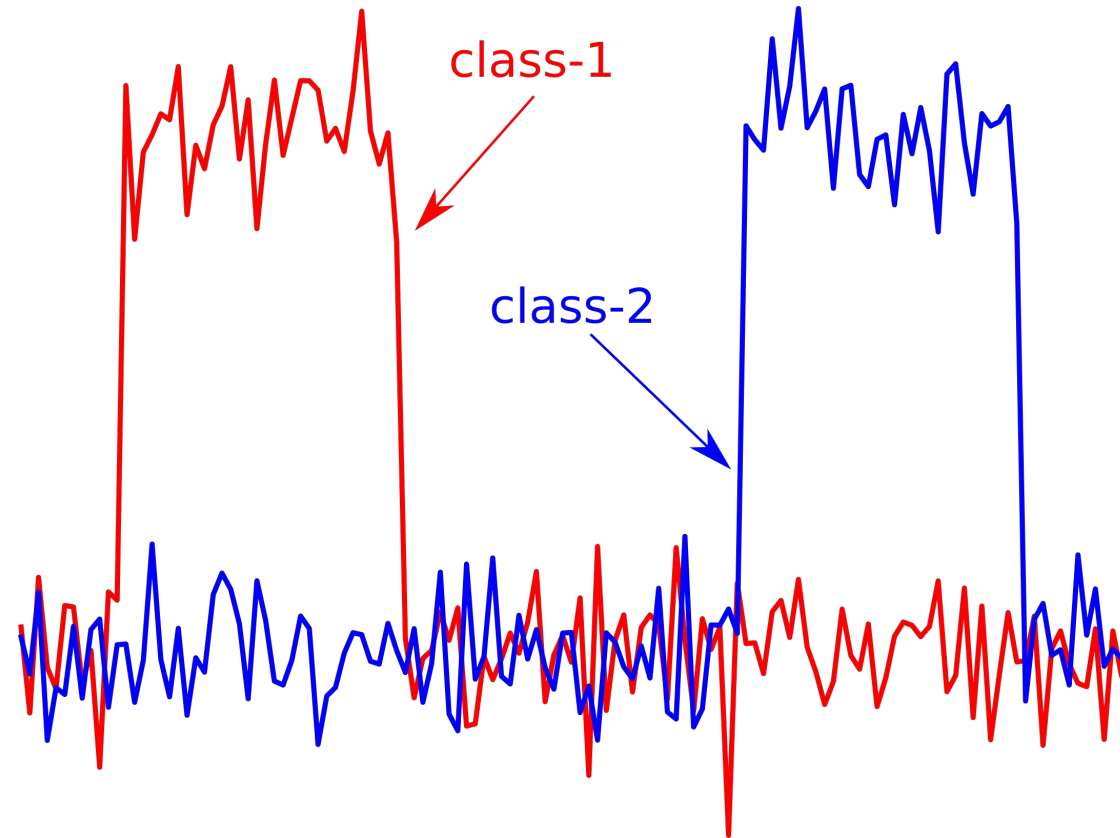
# Studying the size of the ensemble



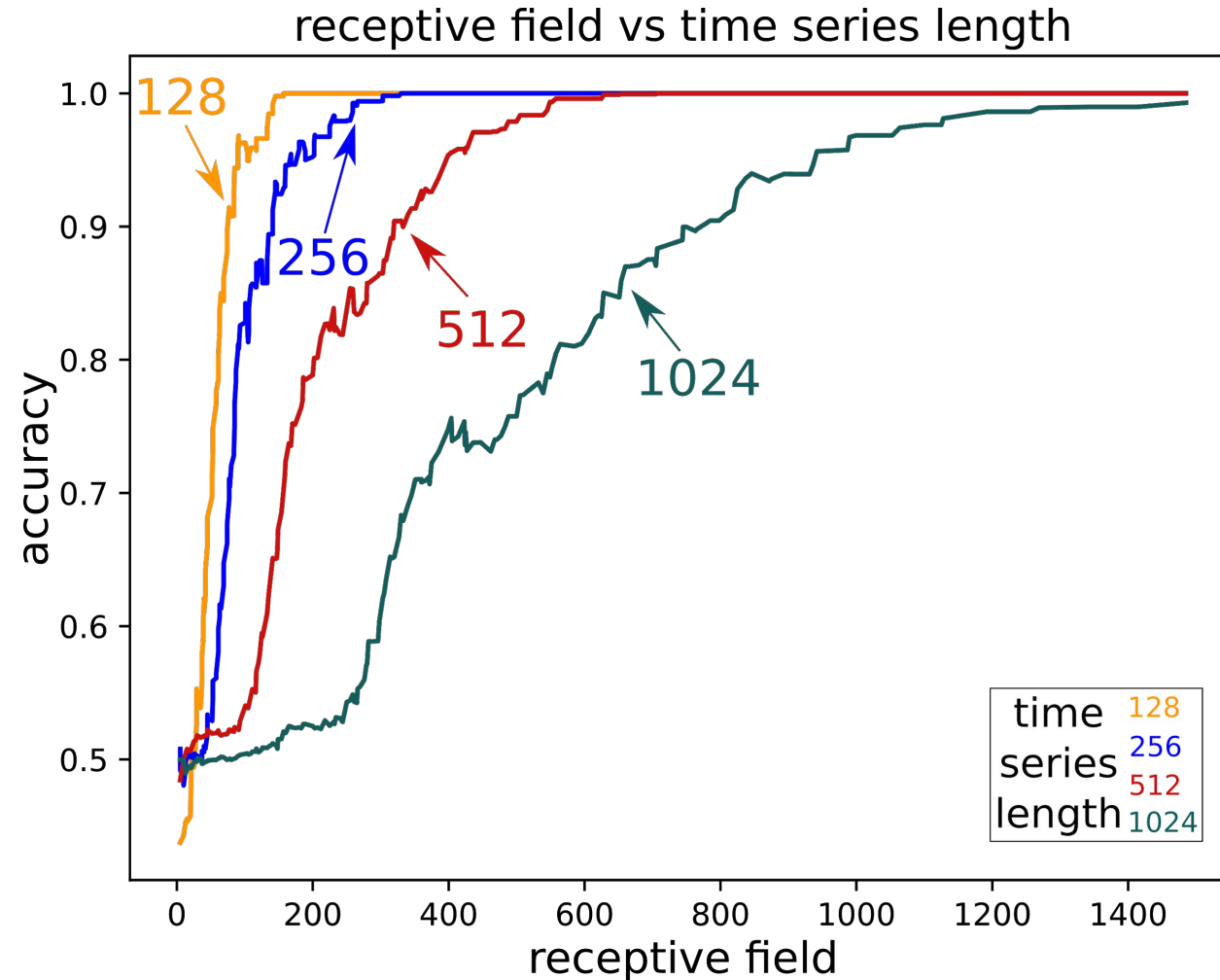Critical difference diagram showing the effect of the number of individual classifiers in InceptionTime

- InceptionTime(*x*) denotes an ensemble of *x* Inception networks
- InceptionTime is equivalent to InceptionTime(5)
- There is no significant improvement for x≥5
  - Again this is due to covariance that start hurting us from 5 elements
  - Therefore we decided to stick with InceptionTime(5)
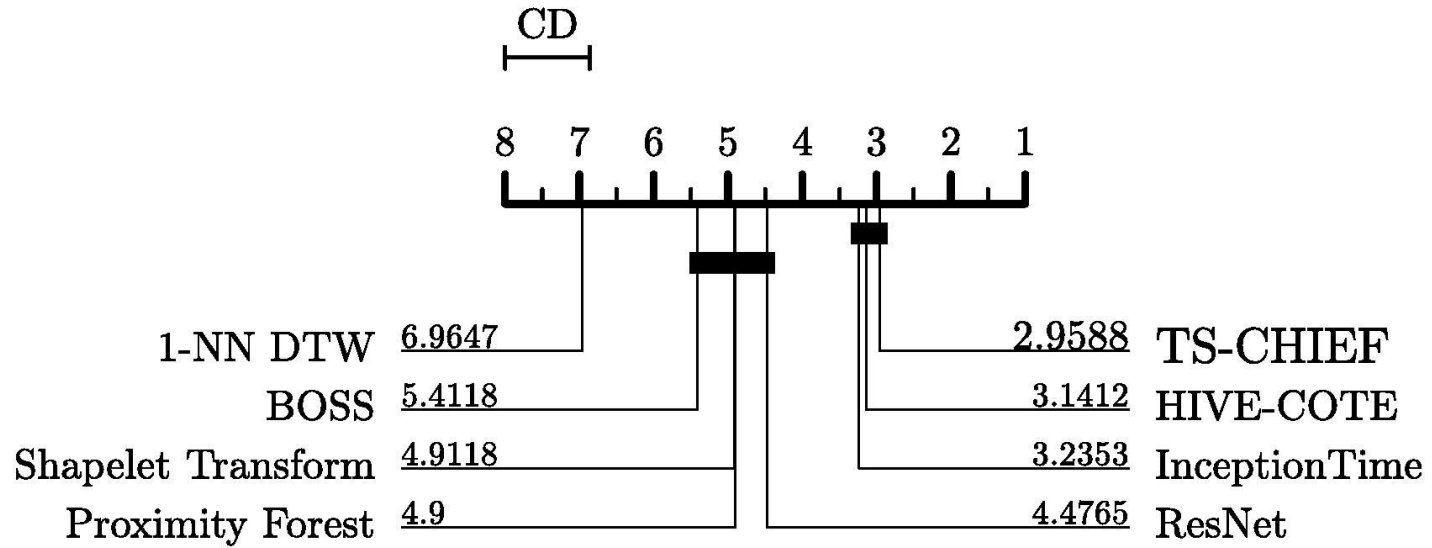
# Hyperparameter study: synthetic dataset



Example of a synthetic binary time series classification problem
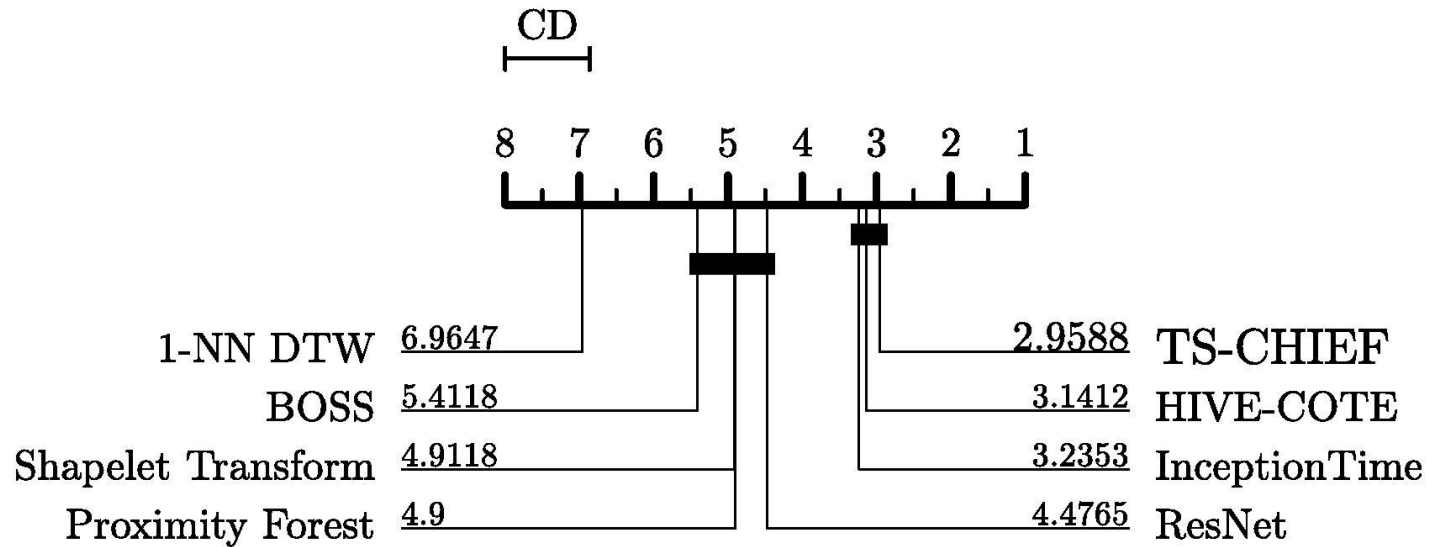
# Hyperparameter study: Receptive Field (RF)



receptive field vs time series length

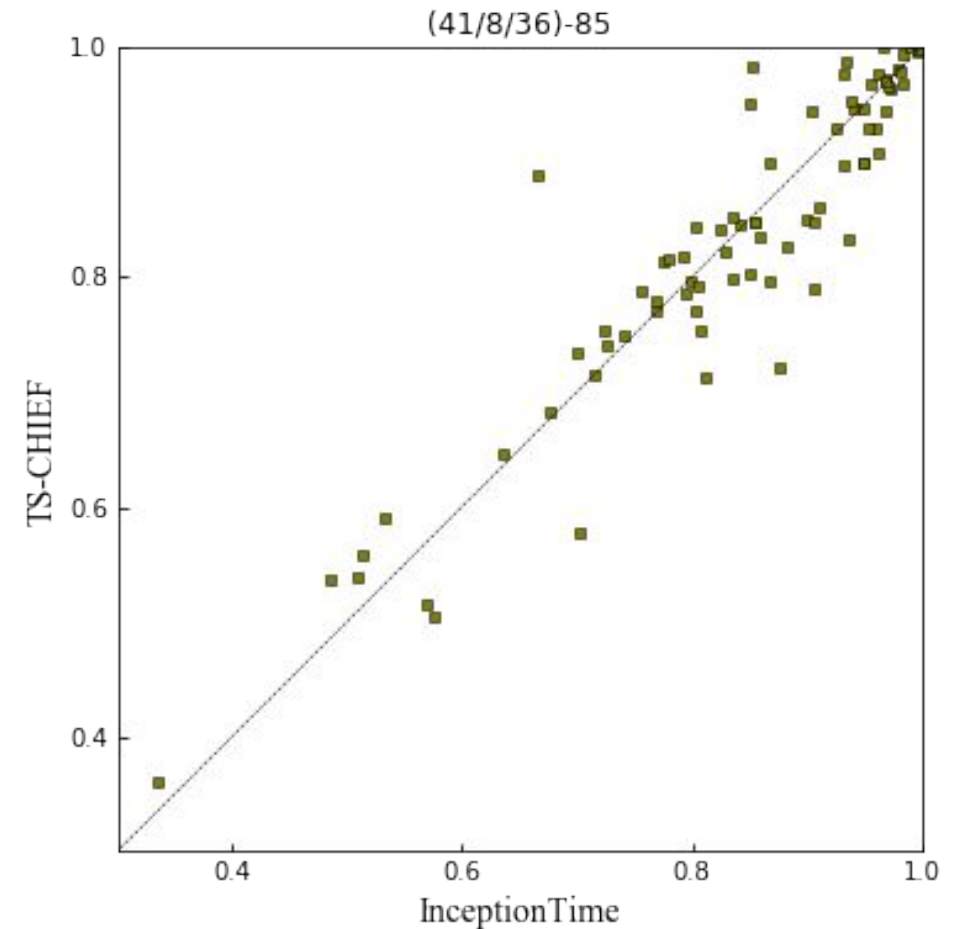A larger receptive field is needed to classify very long time series

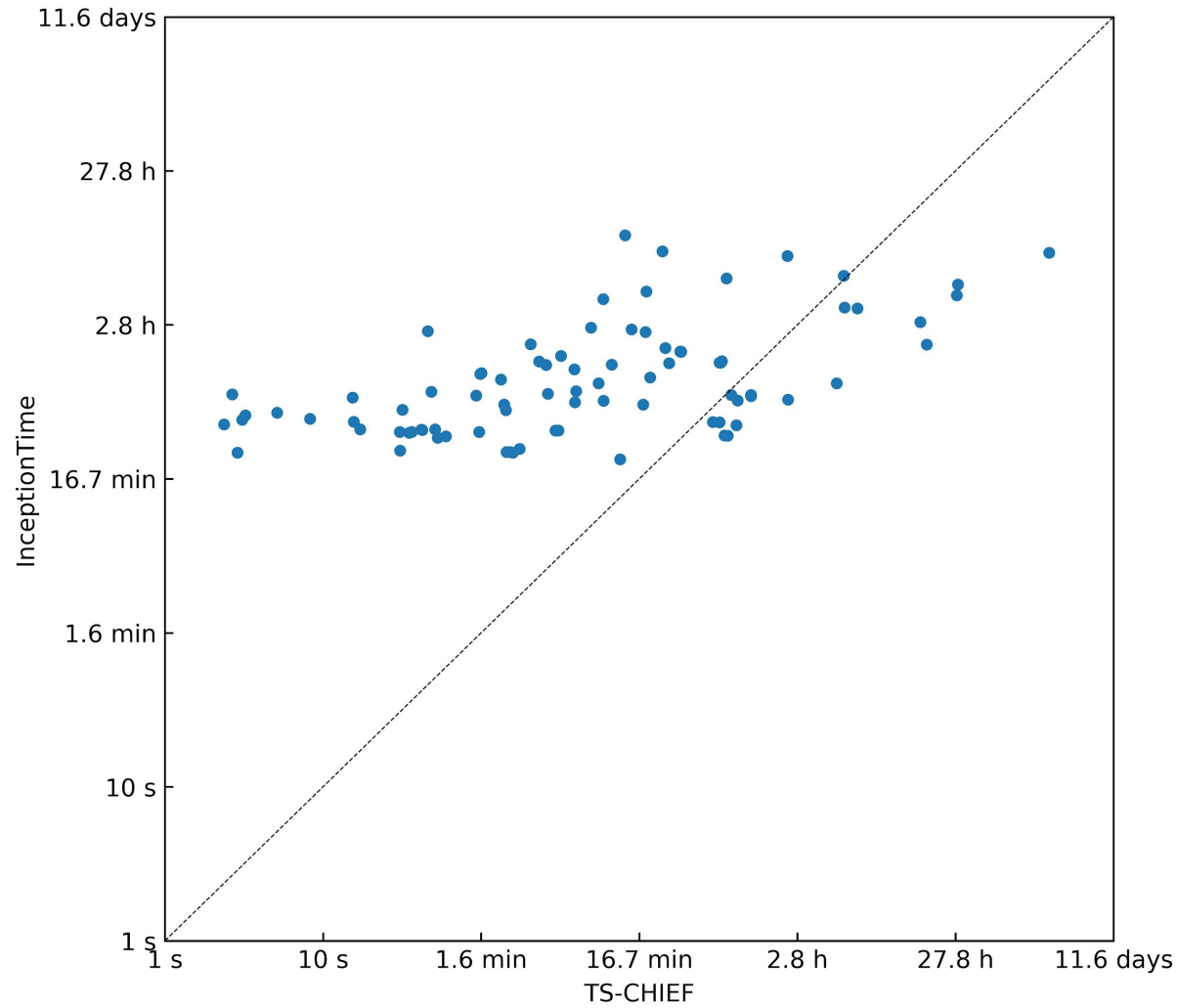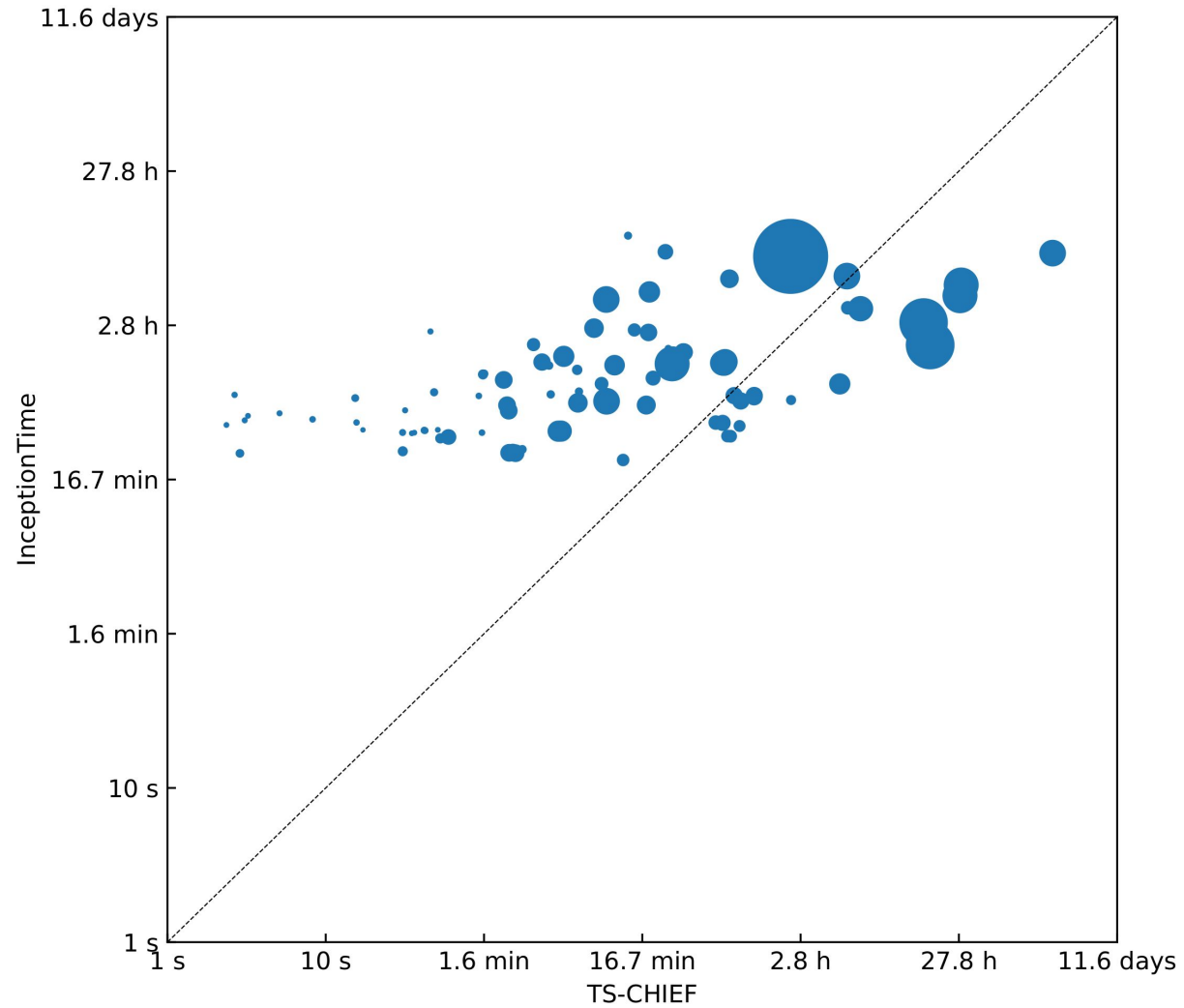# TS-CHIEF vs InceptionTime

# TS-CHIEF vs InceptionTime



- TS-CHIEF wins on average in terms of WDL
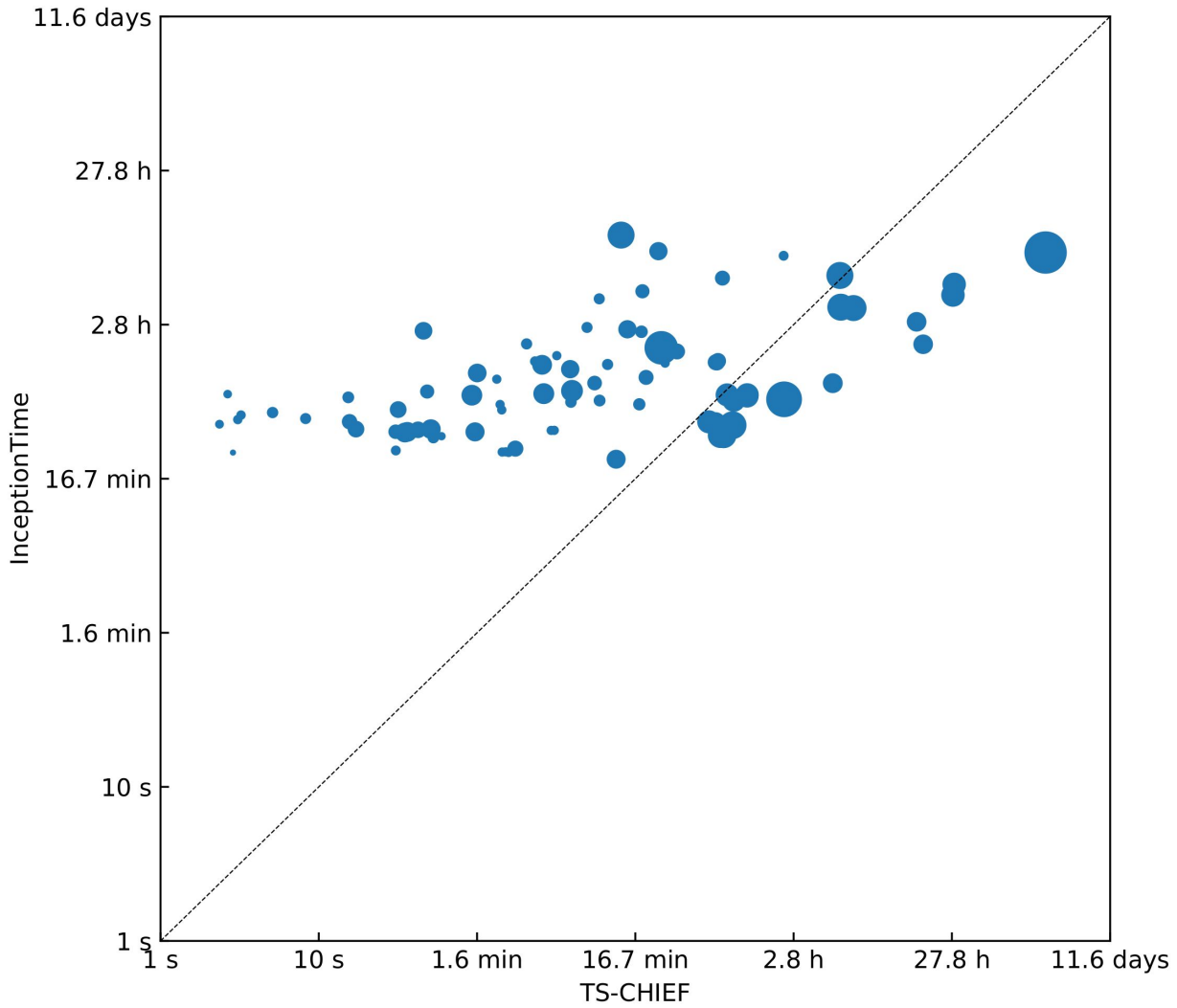- But, lots of big wins with InceptionTime

# TS-CHIEF vs InceptionTime - training time
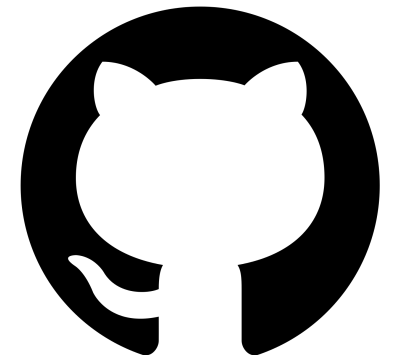
# TS-CHIEF vs InceptionTime - training time

# TS-CHIEF vs InceptionTime - training time

# Conclusions

- Ensemble techniques have revolutionised time series classification with Tony and Jason's group giving us a beacon for research

- TS-Chief combines the efficiencies of tree-based divide-and-conquer with random split selection and the effectiveness of decades' worth of specialised time series techniques

- InceptionTime brings the power and efficiency of deep learning

- Both make state-of-the-art accuracy computationally feasible for large learning tasks

- We believe in reproducible research:

  - Proximity Forest → https://github.com/fpetitjean/ProximityForest

  - TS-CHIEF → https://github.com/dotnet54/TS-CHIEF

  - InceptionTime → https://github.com/hfawaz/InceptionTime

JOIN US IN MELBOURNE!
- 2.5-year postdoc in ML - http://bit.ly/JobsFrancois
- 3 year postdoc+dev in time series

3,600 € per month after-tax

PhD positions available

Send me an email if interested
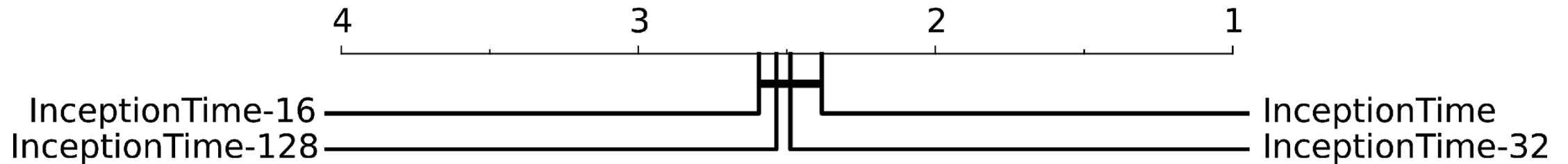→ francois.petitjean@monash.edu

MONASH University

Thank you!

http://francois-petitjean.com

GROUP OF EIGHT AUSTRALIA

# Additional slides
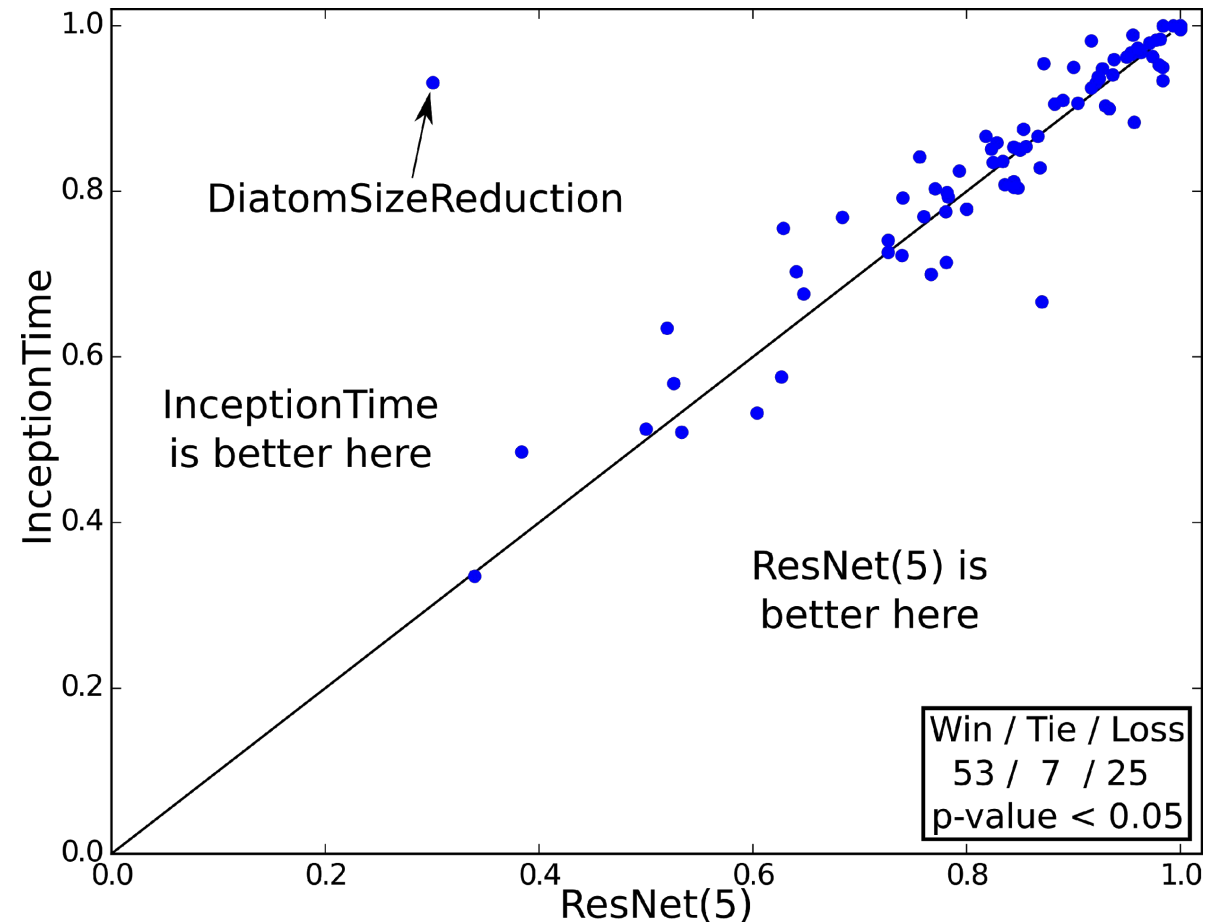
# Hyperparameter study: batch size



Critical difference diagram showing the effect of the batch size hyperparameter value over InceptionTime's average rank

- InceptionTime-*x* denotes InceptionTime with a batch size equal to *x*
- InceptionTime is equivalent to InceptionTime-64 (default value)
- There is no significant difference between the different models
- A value equal to 64 shows a small non-significant superiority
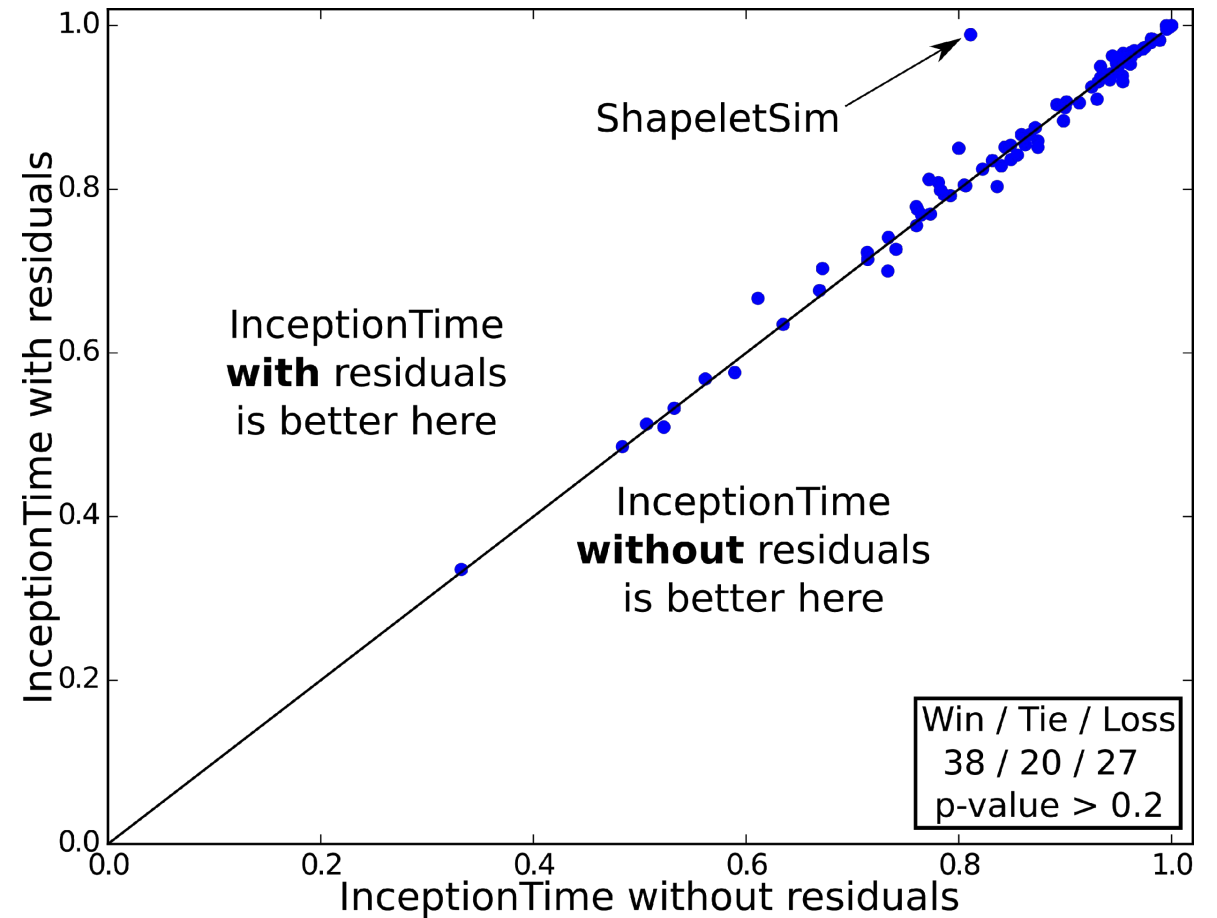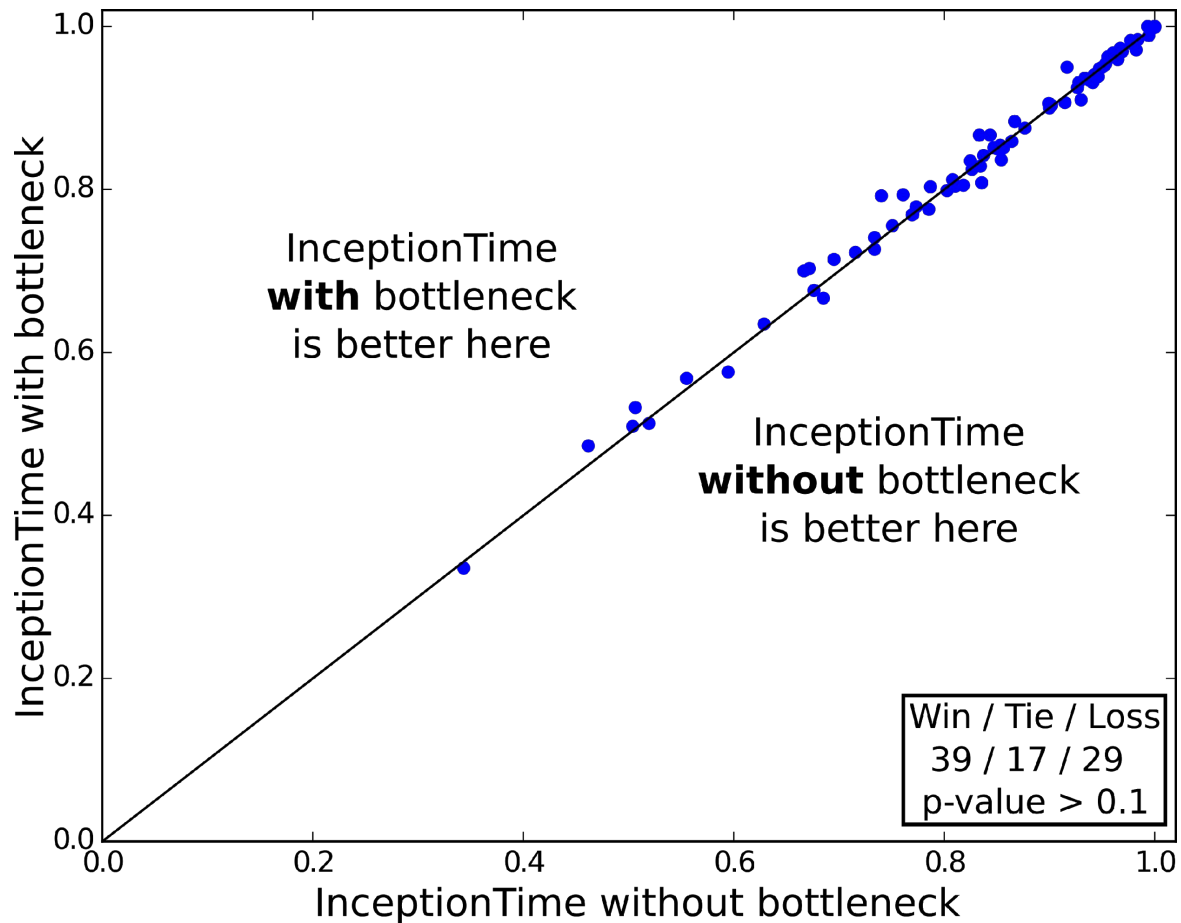- We therefore chose to stick with a batch size equal to 64

# Accuracy plot: InceptionTime vs ResNet(5)

- InceptionTime significantly outperforms ResNet(5) [1]

- For DiatomSizeReduction the main improvement is from using a batch size larger than 1 (which is the case for the ResNet model for this specific dataset)

1. Ismail Fawaz, H., Forestier, G., Weber, J., Idoumghar, L., & Muller, P. (2019). Deep neural network ensembles for time series classification. IEEE International Joint Conference on Neural Networks.
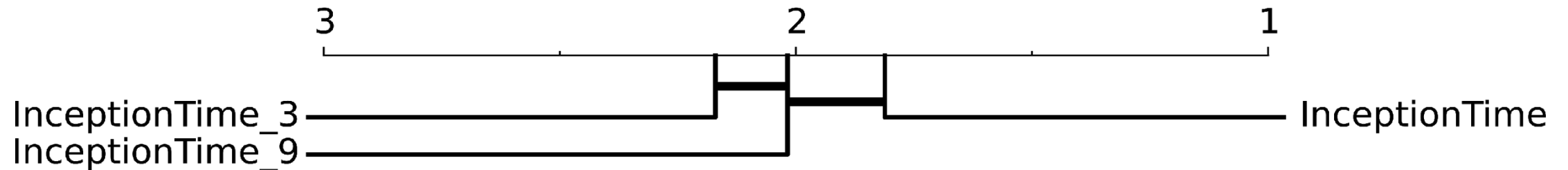
# Hyperparameter study: Bottleneck & residual



Further investigations ShapeletSim indicated that InceptionTime without the residual connections suffered from a severe overfitting.
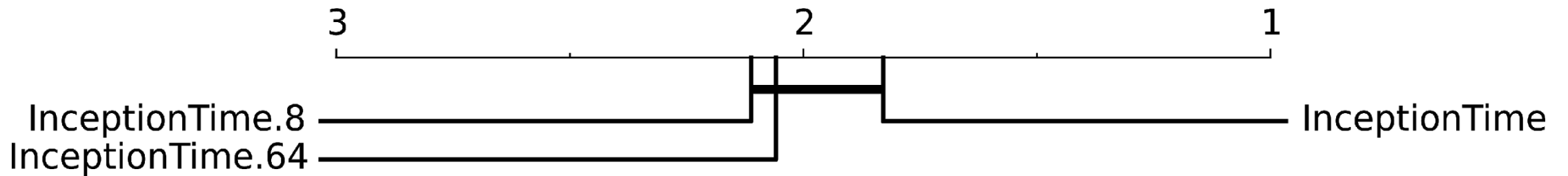
# Hyperparameter study: depth



- InceptionTime_*x* denotes an InceptionTime with *x* layers
- InceptionTime is equivalent to InceptionTime_6 (the default value)
- A shallower model significantly decreases the accuracy
- A deeper model slightly decreases the accuracy
- Therefore we chose to use a network with 6 layers

# Hyperparameter study: number of filters



- InceptionTime:*x* denotes an InceptionTime with *x* filters per module
- InceptionTime is equivalent to a model with 32 filters (default value)
- More filters showed a significant decrease in accuracy
- Less filters showed a slight decrease in accuracy
- This hyperparameter affects significantly the complexity of the model

# Hyperparameter study: filter length



- InceptionTime.*x* denotes a model with a filter length equal to *x*

- InceptionTime is equivalent to a model with a filter length equal to 32

- The default value (32) showed a slight advantage

- Although larger values will produce a larger RF, these experiments showed that this hyperparameter should be carefully chosen

# Receptive Field (RF) of a neural network

$$1 + \sum_{i=1}^{d}(k_i - 1)$$

- $d$ represents the depth of the network
- $k_i$ represents the length of the filters in $i^{th}$ layer
- The stride is considered to be equal to 1
- RF can then be increased by either controlling $d$ or $k_i$

For images, a large RF is needed to capture more context [1]

1. Luo, W., Li, Y., Urtasun, R., & Zemel, R. (2016). Understanding the effective receptive field in deep convolutional neural networks. In *Advances in neural information processing systems*.