

# BEFORE DOING EXPERIMENTS

Qi Chen

[Qi.Chen@ecs.vuw.ac.nz](mailto:Qi.Chen@ecs.vuw.ac.nz)

Bing Xue

[Bing.Xue@ecs.vuw.ac.nz](mailto:Bing.Xue@ecs.vuw.ac.nz)

Harith Al-Sahaf

[Harith.Al-Sahaf@ecs.vuw.ac.nz](mailto:Harith.Al-Sahaf@ecs.vuw.ac.nz)



# Outline

- Experiment Settings
- Parameter Settings
- Results
- Codes and Programs
- ~~Grid Computing~~



# Experiment Settings

- Baseline and Benchmark methods: select appropriate models for comparison
  - Task relevant: feature learning, classification, regression, clustering, ...
  - **Baseline** methods: simple, well-established methods, provide a reference point
  - **Original** method: to demonstrate the improvement over the predecessor
  - **State-of-the-art** methods: currently achieve the best performance on a specific task, also consider
    - ✓ time of publication: more likely to be more recent publications, a general rule of thumb, 2-3 years old can be considered "old"
    - ✓ publication venue, top-tier journal and conference



# Experiment Settings

- Benchmark Datasets: use datasets commonly used for benchmarking, number of datasets: generally more than 10
- Training and Test Sets
  - Which method to create training and test sets
    - Hold-out/training-test split: with a splitting ratio of 80:20 or 70:30
    - K-fold cross-validation (sampling without replacement)
    - Bootstrap (sampling with replacement)
  - Splitting seed
  - Stratified splitting to maintain class-ratio



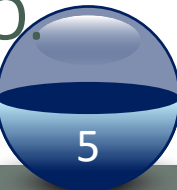
# Experiment Settings (continue)

- Number of runs:
  - at least **30** runs (using 30 different seeds) => WHY?
  - or 50 runs unless your experiment is too time-consuming.
- Performance evaluation:
  - Measure: error rate, accuracy.
  - What to report if you have (10-fold cross-validation x 30 runs) results?



# Random Seeds For EC methods

- DO record your random seeds to make sure you can **re-produce** the same results later if needed
  - Do NOT use clock time as the random seed
- Use the same random seeds to compare two different versions of the same approach:
  - E.g. two different GP algorithms: GP1 and GP2, run both of them for 50 times. Please make sure you use the same 50 random seeds for GP1 and GP2 to let them have the same starting points for fair comparisons.
  - It is not necessary to use the same random seeds if you compare GP with PSO.



# Parameter Settings

- Start with commonly used parameter settings from the literature, or settings recommended by good papers.
  - Do not randomly pick up some parameters values unless you have good reasons to use them.
- Figure out what the parameter values exactly mean
- Parameter tuning: one aspect at a time



# Results

- Please **record all** useful results
  - Eg, the *gbest* in PSO, the best program from GP, the training, testing performances in each run (you may further check with your supervisors)
  - Evolutionary plots on the training and test sets: performance of the best individuals at every generation
  - **Computational (training) time of each run**: first generation to the last generation
    - not include test process
  - Model Size, #feature ..
  - Evolved Models



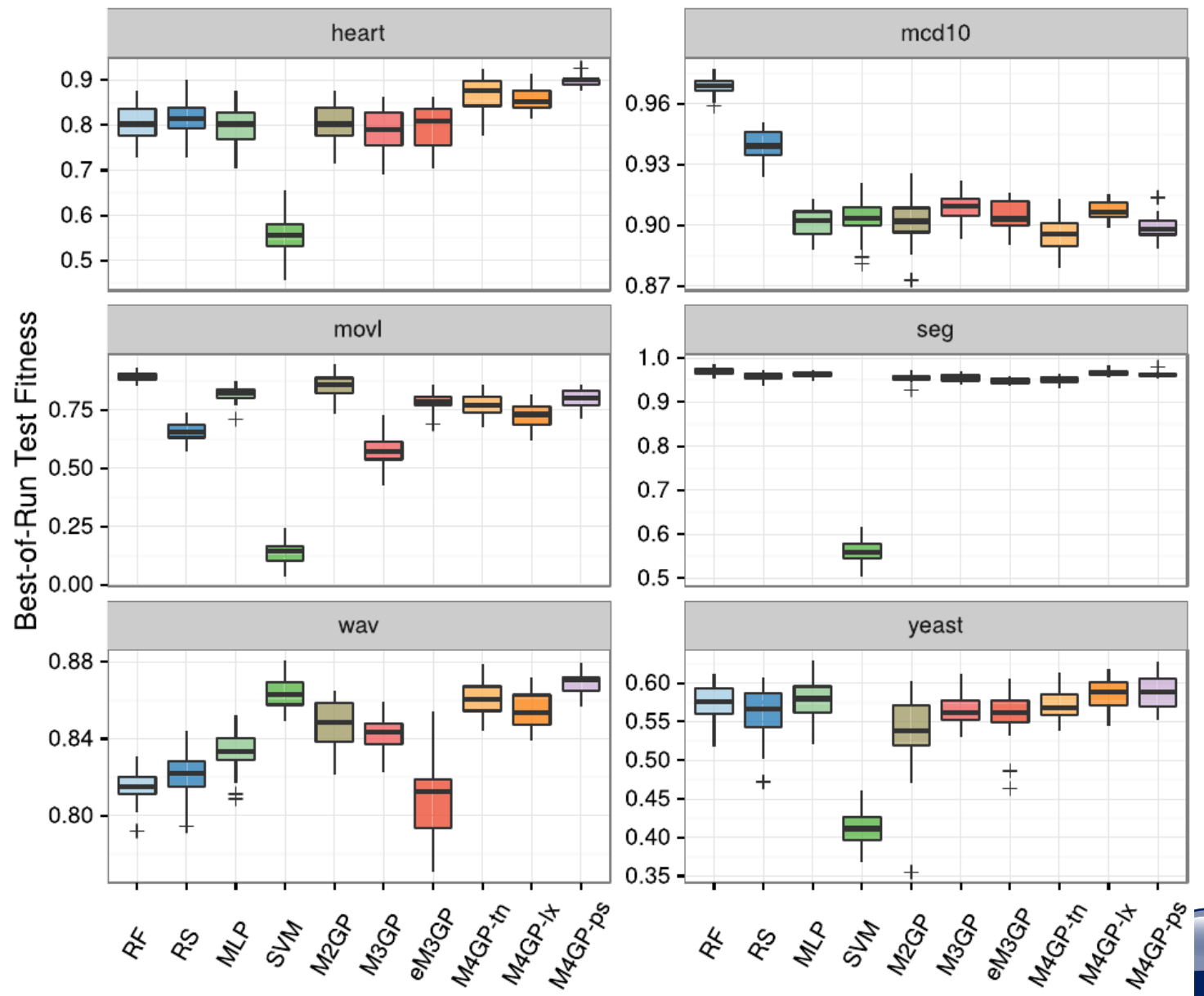
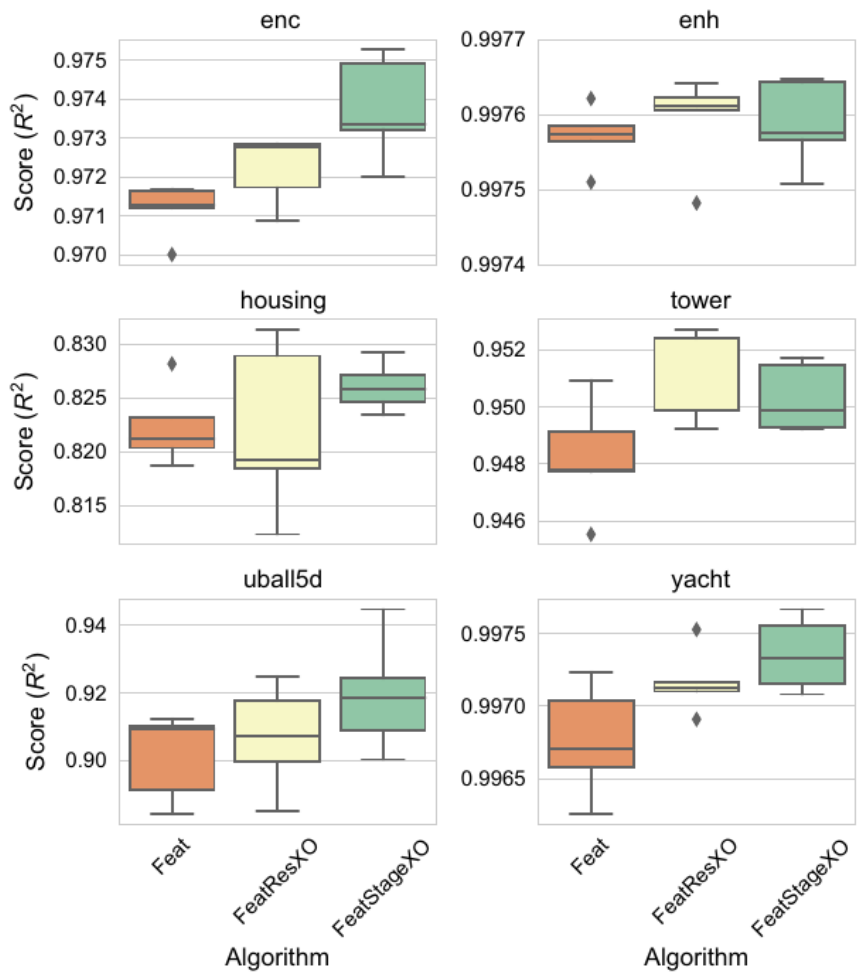


**Table 3**

Comparison of best-of-run median test accuracy for the benchmark problems. The best result is highlighted in bold. Significant ( $p < 0.01$  according to a pairwise Wilcoxon rank-sum test with Holm correction) improvements with respect to each method is denoted by  $a - j$  according to the method labels.

Method	Heart	IM-3	IM-10	Movl	Seg	Vowel	Wav	Yeast
<sup>a</sup> RF	<sup>d</sup> 80.2	94.8	<sup>bcd</sup> <sub>efghij</sub> <b>96.9</b>	<sup>bcd</sup> <sub>efghij</sub> <b>89.4</b>	<sup>bcd</sup> <sub>efghij</sub> <b>97.3</b>	<sup>bcd</sup> <sub>g</sub> 89.4	81.5	<sup>d</sup> <sub>e</sub> 57.5
<sup>b</sup> RS	<sup>d</sup> 81.5	92.8	<sup>cde</sup> <sub>efghij</sub> 93.9	<sup>d</sup> <sub>f</sub> 65.7	<sup>d</sup> <sub>gj</sub> 96.0	<sub>g</sub> 82.8	82.2	<sup>d</sup> 56.6
<sup>c</sup> MLP	<sup>d</sup> 80.2	95.9	90.2	<sup>bd</sup> <sub>efghij</sub> 82.5	<sup>d</sup> <sub>gj</sub> 96.3	<sub>g</sub> 82.5	<sup>ab</sup> <sub>g</sub> 83.3	<sup>d</sup> <sub>e</sub> 58.0
<sup>d</sup> SVM	55.6	93.8	90.4	14.4	55.8	81.8	<sup>abc</sup> <sub>efgh</sub> 86.3	41.1
<sup>e</sup> M2GP	<sup>d</sup> 80.2	93.8	90.2	<sup>bcd</sup> <sub>efghij</sub> 85.9	<sup>d</sup> <sub>g</sub> 95.6	<sup>cd</sup> <sub>g</sub> 85.9	<sup>abc</sup> <sub>g</sub> 84.9	<sup>d</sup> 53.8
<sup>f</sup> M3GP	<sup>d</sup> 79.0	95.4	<sup>c</sup> <sub>ij</sub> 91.0	<sup>d</sup> 57.1	<sup>d</sup> 95.6	<sup>abcde</sup> <sub>g</sub> 93.8	<sup>ab</sup> <sub>g</sub> 84.3	<sup>d</sup> 56.2
<sup>g</sup> eM3GP	<sup>d</sup> 80.9	93.3	<sub>j</sub> 90.3	<sup>bd</sup> <sub>fh</sub> 78.6	<sup>d</sup> 94.7	78.6	81.2	<sup>d</sup> 56.2
<sup>h</sup> M4GP-lx	<sup>abcde</sup> <sub>fg</sub> 85.2	<sup>abcde</sup> <sub>fg</sub> <b>97.9</b>	<sub>ij</sub> 90.7	<sup>bd</sup> <sub>f</sub> 73.1	<sup>bde</sup> <sub>fji</sub> 96.6	<sup>abcde</sup> <sub>fg</sub> 95.6	<sup>abc</sup> <sub>fg</sub> 85.3	<sup>d</sup> <sub>e</sub> <b>58.9</b>
<sup>i</sup> M4GP-ps	<sup>abcde</sup> <sub>efghij</sub> <b>90.1</b>	<sup>abcde</sup> <sub>fg</sub> 97.9	89.8	<sup>bd</sup> <sub>fh</sub> 80.1	<sup>d</sup> <sub>gj</sub> 96.1	<sup>abcde</sup> <sub>efghij</sub> <b>97.5</b>	<sup>abc</sup> <sub>efghij</sub> <b>87.1</b>	<sup>d</sup> <sub>e</sub> <sub>fg</sub> 58.9
<sup>j</sup> M4GP-tn	<sup>abcde</sup> <sub>fg</sub> 87.7	<sup>abcde</sup> <sub>fg</sub> 97.9	89.6	<sup>bd</sup> <sub>fh</sub> 76.9	<sup>d</sup> 95.1	<sup>abcde</sup> <sub>fg</sub> 96.0	<sup>abc</sup> <sub>efgh</sub> 86.0	<sup>d</sup> 56.8





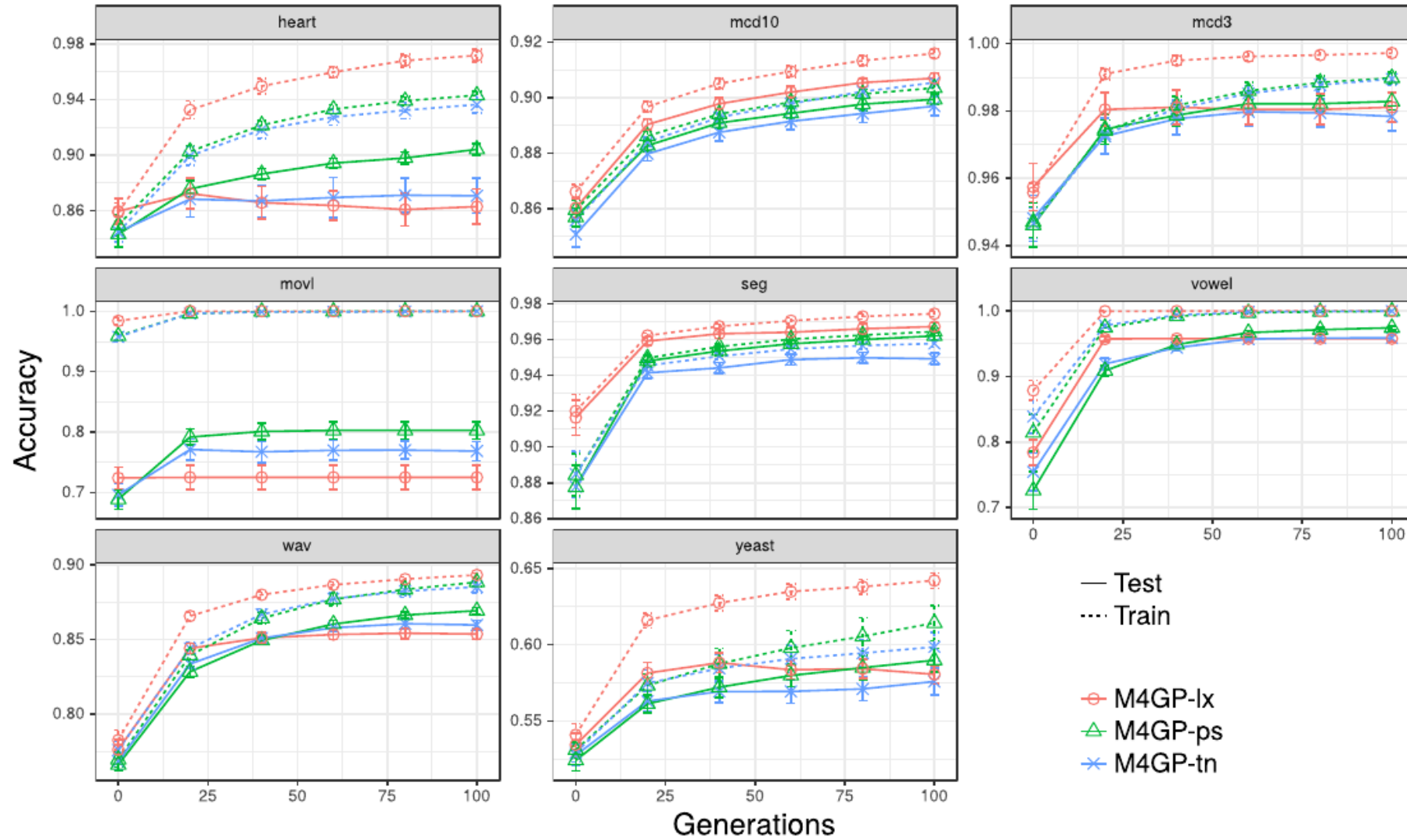
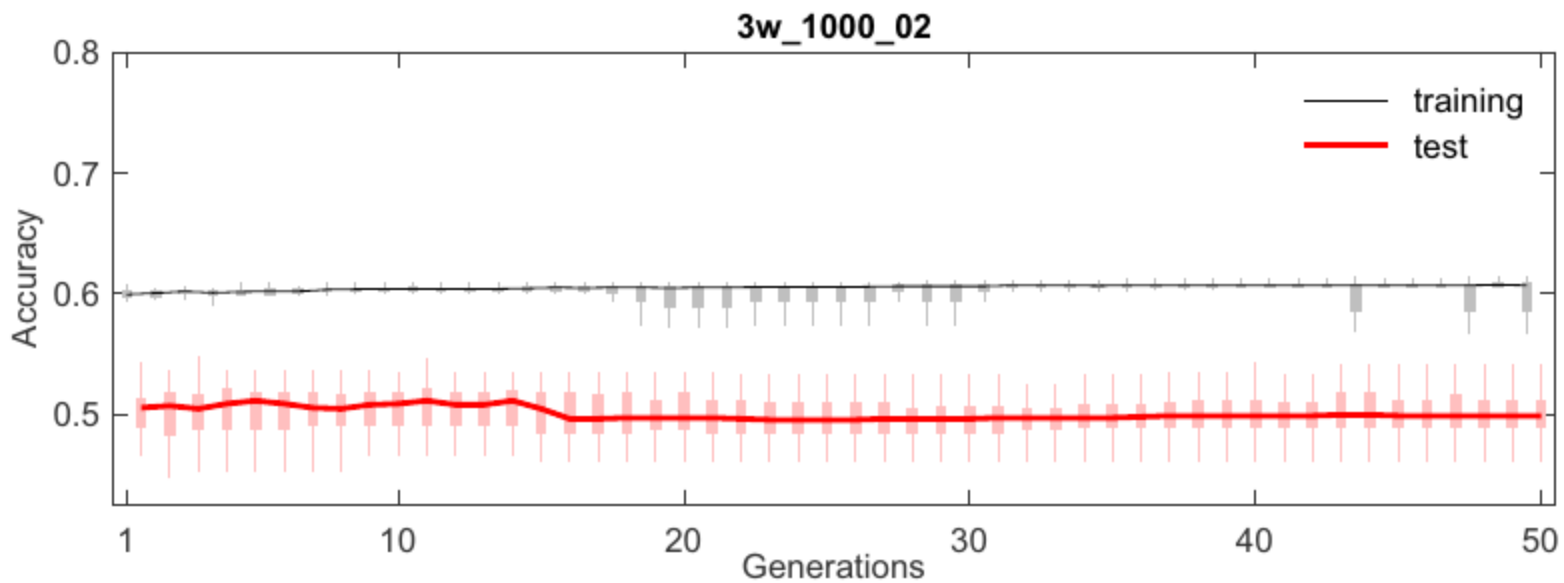
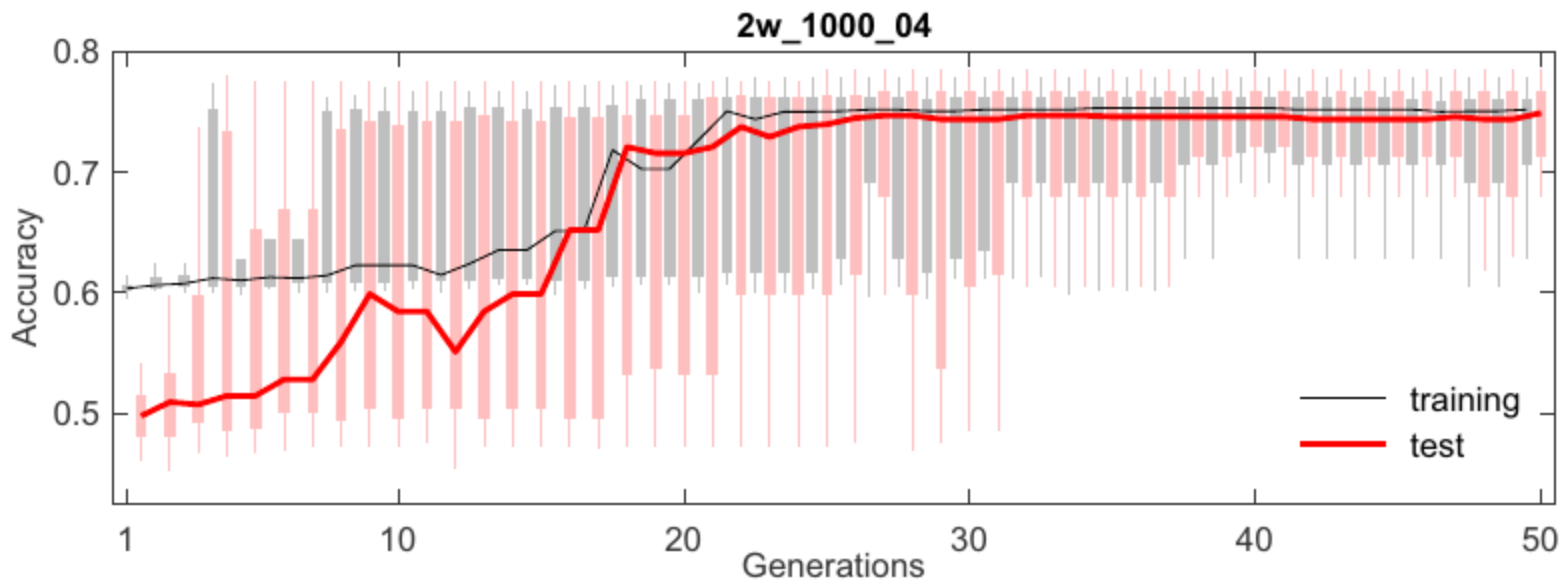
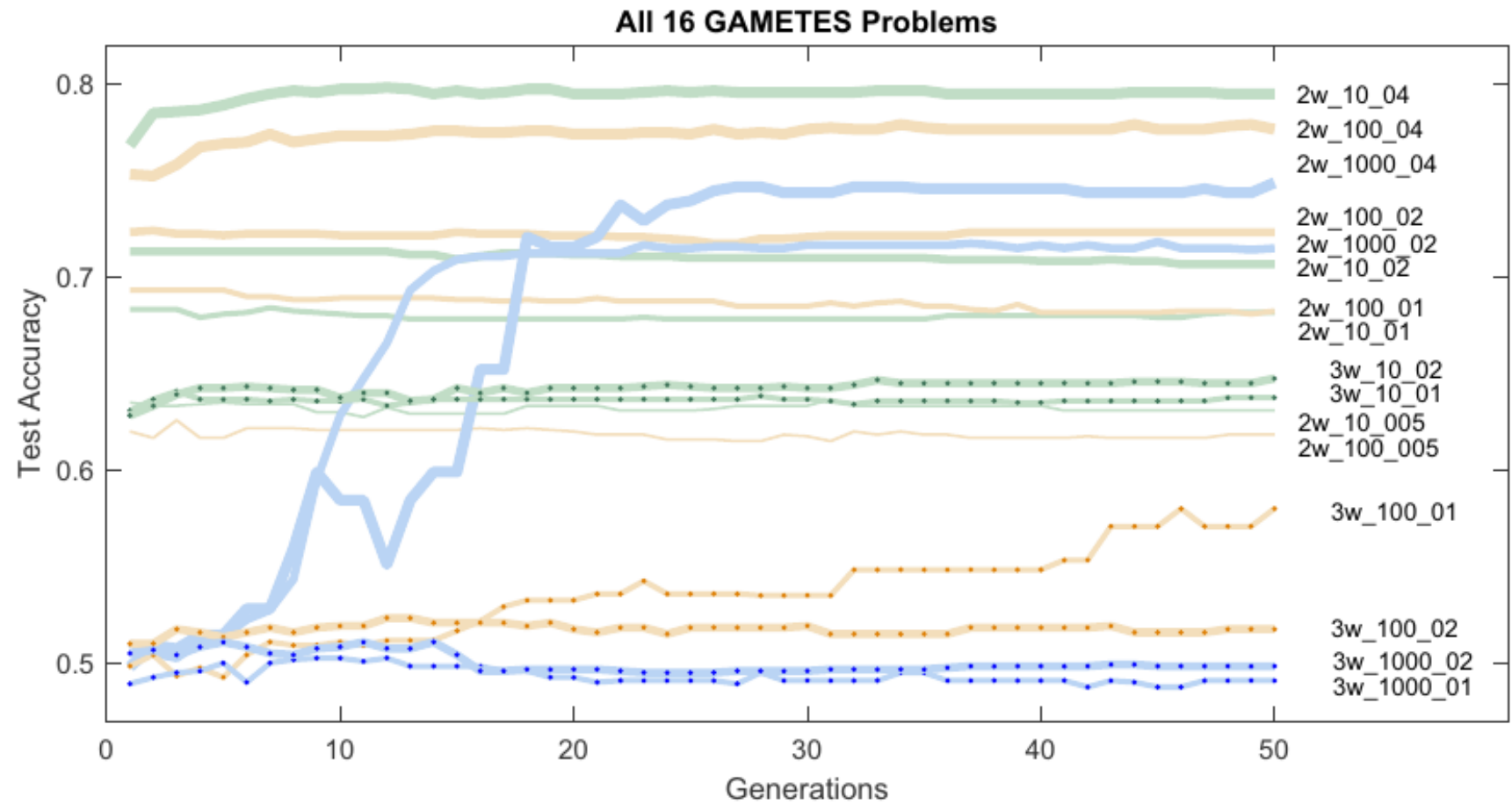


Fig. 5. Mean convergence characteristics of M4GP on the training set (dotted lines) and test set (solid lines) over 100 generations. Shapes indicate different selection methods. Error bars denote the confidence intervals.

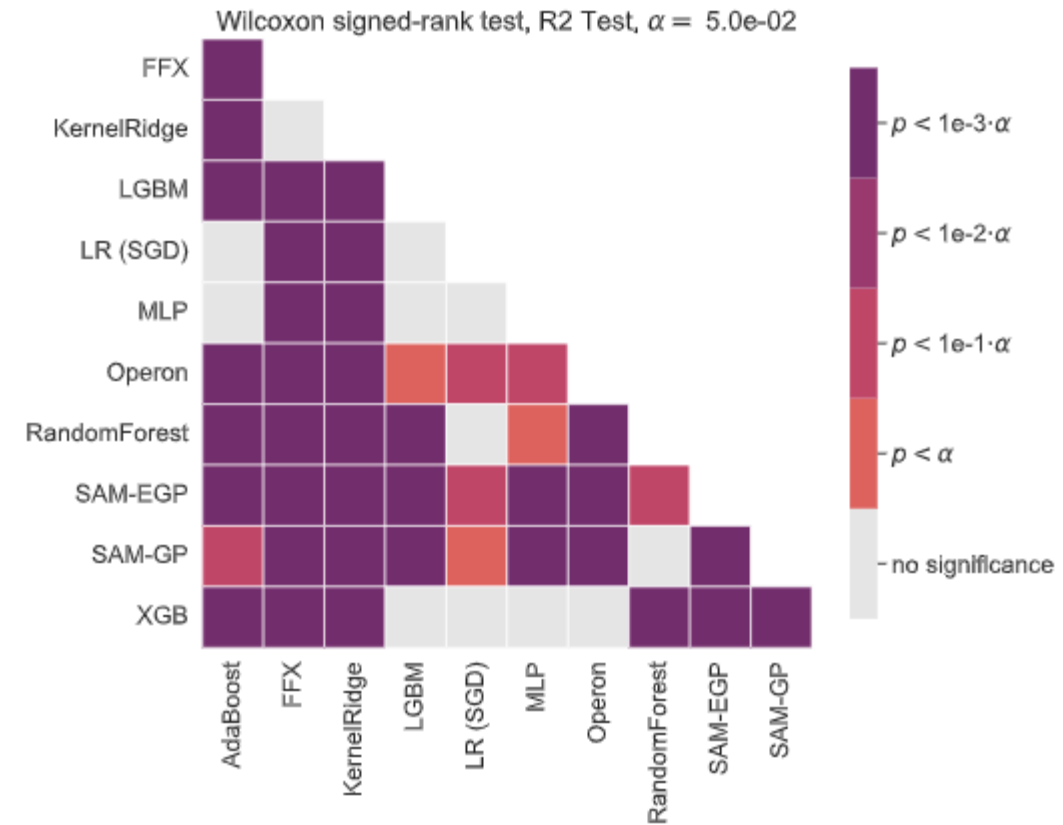


**Fig. 8** Evolution of test accuracy on all the GAMETES problems. Lines are green/yellow/blue for 10/100/1000-feature problems; lines are thicker for higher signal-to-noise ratios; lines are dotted for the three-way problems. The names of the problems on the right appear by the same order as the lines



# Statistical significance test results

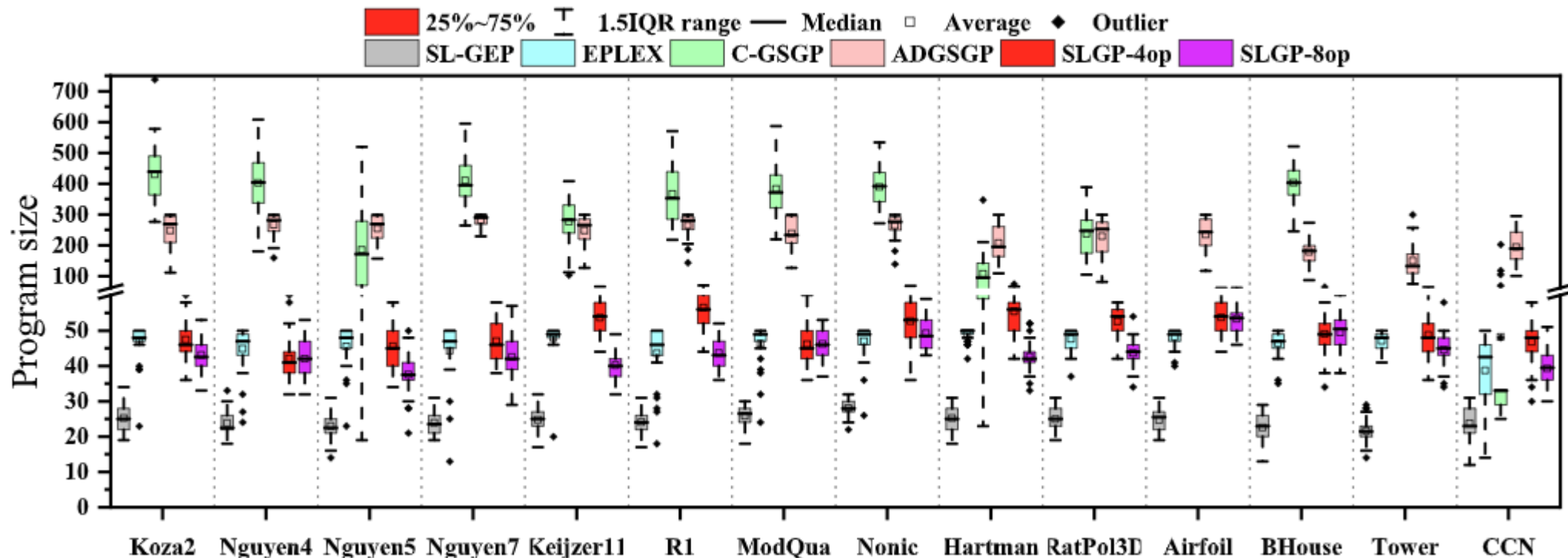
	SGP_E	IGP	IGP_E	FIGP	FIGP_E	SGP_E	IGP	IGP_E	FIGP	FIGP_E
	547_no2					1193_BNG_lowbwt				
SGP	3.08e-1	9.03e-1	3.71e-1	<b>4.60e-4</b>	<b>8.86e-5</b>	7.15e-1	<b>1.85e-2</b>	3.93e-1	3.81e-1	<b>3.73e-9</b>
SGP_E		3.18e-1	6.70e-2	<b>4.97e-5</b>	<b>1.22e-5</b>		<b>2.09e-2</b>	7.61e-1	1.98e-1	<b>8.01e-8</b>
IGP			2.29e-1	<b>1.37e-2</b>	<b>5.01e-3</b>			<b>2.02e-3</b>	1.98e-1	<b>1.19e-6</b>
IGP_E				<b>1.64e-2</b>	<b>3.74e-3</b>				<b>2.62e-2</b>	<b>9.31e-9</b>
FIGP					3.49e-1					<b>1.30e-8</b>
	503_wind					195_auto_price				
SGP	2.89e-1	<b>4.18e-4</b>	6.85e-1	<b>7.61e-3</b>	<b>3.85e-7</b>	<b>5.59e-9</b>	<b>1.72e-3</b>	3.09e-1	<b>9.31e-9</b>	<b>1.99e-6</b>
SGP_E		<b>2.37e-3</b>	1.91e-1	<b>3.64e-2</b>	<b>1.86e-9</b>		<b>2.05e-7</b>	<b>1.68e-6</b>	<b>1.86e-9</b>	<b>1.86e-9</b>
IGP			<b>7.98e-4</b>	8.24e-1	<b>1.37e-4</b>			<b>2.56e-3</b>	<b>1.30e-8</b>	<b>8.33e-7</b>
IGP_E				<b>4.60e-4</b>	<b>1.86e-9</b>				<b>1.30e-8</b>	<b>2.32e-4</b>
FIGP					<b>4.67e-3</b>					<b>3.86e-7</b>
	583_fri_c1_1000_50					588_fri_c4_1000_100				
SGP	1.84e-1	6.36e-2	9.19e-1	<b>3.22e-3</b>	<b>4.18e-4</b>	<b>2.62e-2</b>	2.53e-1	4.40e-1	<b>2.02e-3</b>	3.09e-1
SGP_E		<b>4.84e-4</b>	<b>3.64e-2</b>	<b>8.86e-5</b>	<b>2.76e-6</b>		<b>2.83e-4</b>	1.29e-1	<b>1.68e-6</b>	<b>2.62e-4</b>
IGP			1.14e-1	2.13e-1	<b>2.37e-3</b>			<b>2.34e-2</b>	<b>2.62e-2</b>	5.03e-1
IGP_E				<b>5.01e-3</b>	<b>1.40e-5</b>				<b>1.60e-5</b>	9.20e-2
FIGP					2.21e-1					<b>9.30e-3</b>
	4544_GeographicalOriginalofMusic					505_tecator				
SGP	1.24e-1	2.67e-1	1.29e-1	<b>3.73e-9</b>	<b>1.70e-4</b>	<b>1.30e-7</b>	<b>1.30e-2</b>	9.61e-2	<b>8.00e-3</b>	<b>9.31e-9</b>
SGP_E		6.70e-1	9.68e-1	<b>1.86e-8</b>	5.01e-3		3.82e-1	<b>6.20e-3</b>	1.19e-1	5.77e-2
IGP			8.71e-1	<b>3.73e-9</b>	<b>1.72e-3</b>			3.39e-1	7.77e-1	5.49e-2
IGP_E				<b>1.86e-9</b>	<b>1.11e-4</b>				5.43e-1	<b>1.70e-4</b>
FIGP					<b>5.15e-6</b>					<b>1.21e-2</b>



# Model Size, #feature ...

AVERAGE FEATURE SUBSET SIZE OBTAINED BY VGS-MOEA AND ITS TWO VARIANTS

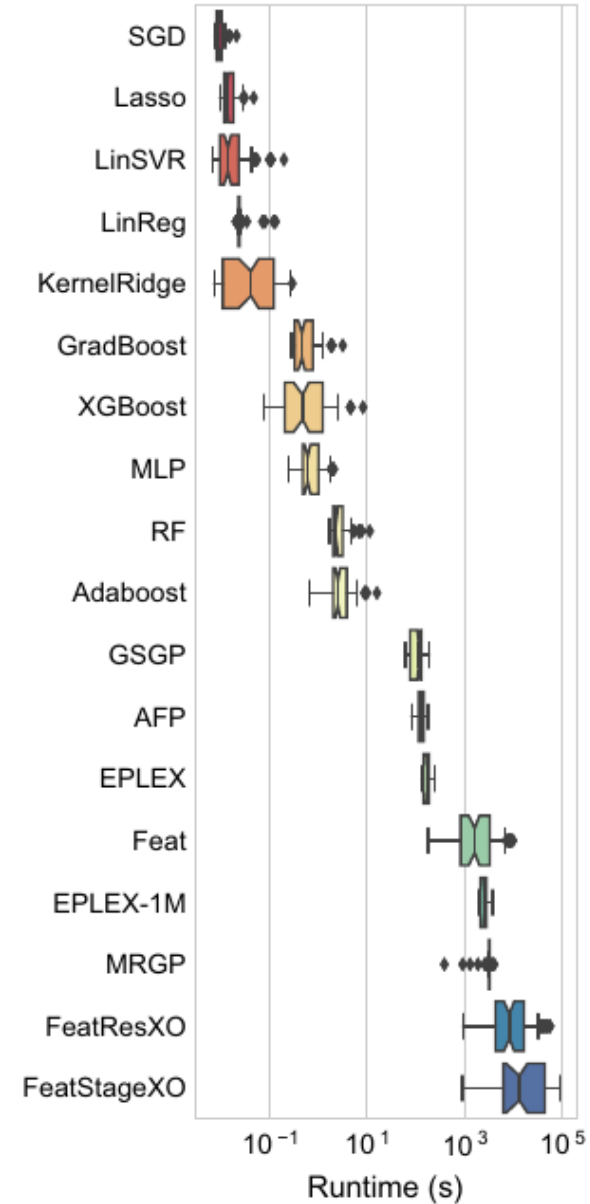
Data Set	VGS-MOEA-NAR	VGS-MOEA-NGO	VGS-MOEA
COIL20	94.40 ± 36.32(+)	160.20 ± 44.34(+)	<b>59.03 ± 25.91</b>
SRBCT	155.13 ± 25.56(+)	200.07 ± 41.58(+)	<b>29.33 ± 21.57</b>
PCMAC	368.13 ± 138.28(+)	343.23 ± 109.86(+)	<b>50.33 ± 10.67</b>
lymphoma	326.33 ± 70.82(+)	407.10 ± 70.59(+)	<b>70.70 ± 47.53</b>
GLIOMA	412.80 ± 100.16(+)	558.87 ± 117.29(+)	<b>23.97 ± 26.06</b>
BASEHOCK	465.20 ± 151.87(+)	475.87 ± 153.12(+)	<b>65.57 ± 19.56</b>
TOX_171	570.53 ± 166.54(+)	776.63 ± 208.77(+)	<b>102.07 ± 72.61</b>
Brain1	495.63 ± 204.95(+)	685.07 ± 138.21(+)	<b>36.47 ± 31.56</b>
leukemia	436.30 ± 68.84(+)	552.97 ± 132.53(+)	<b>24.17 ± 18.69</b>



# Computational Time

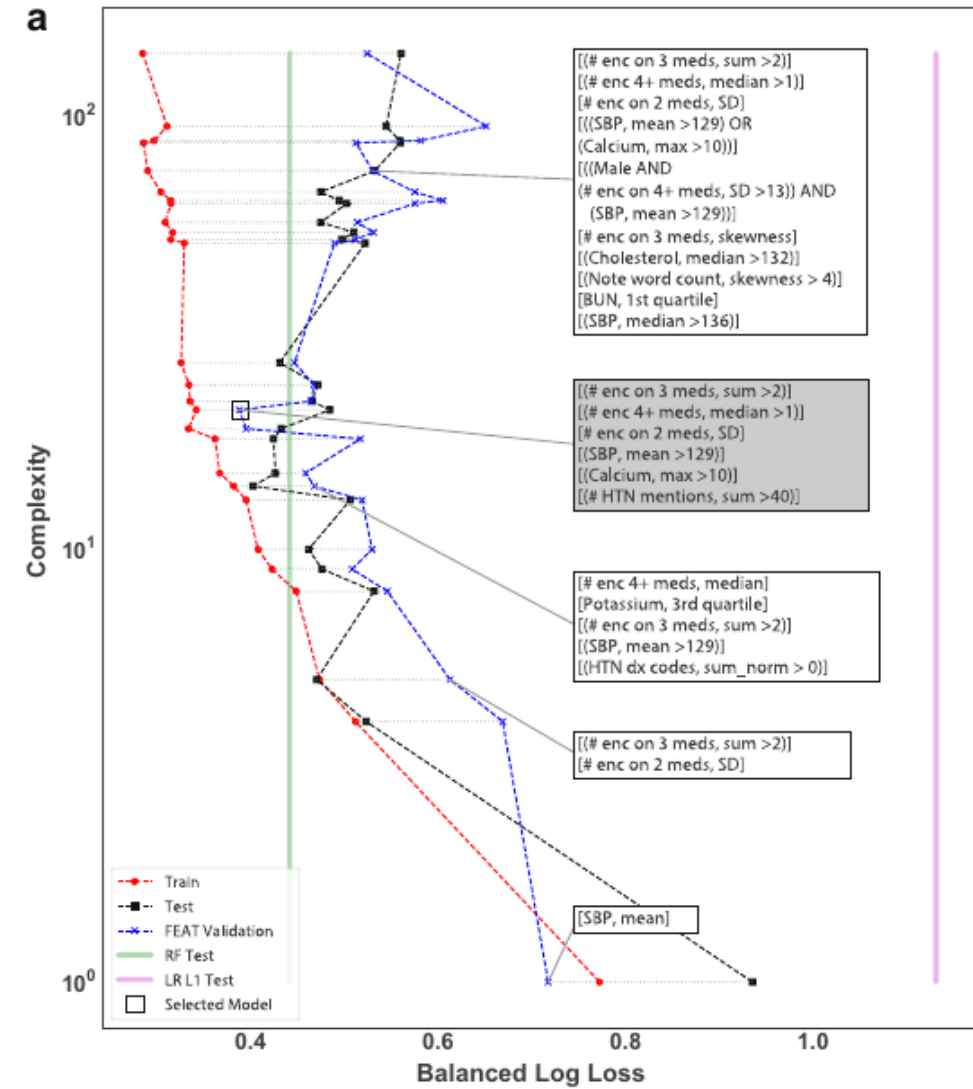
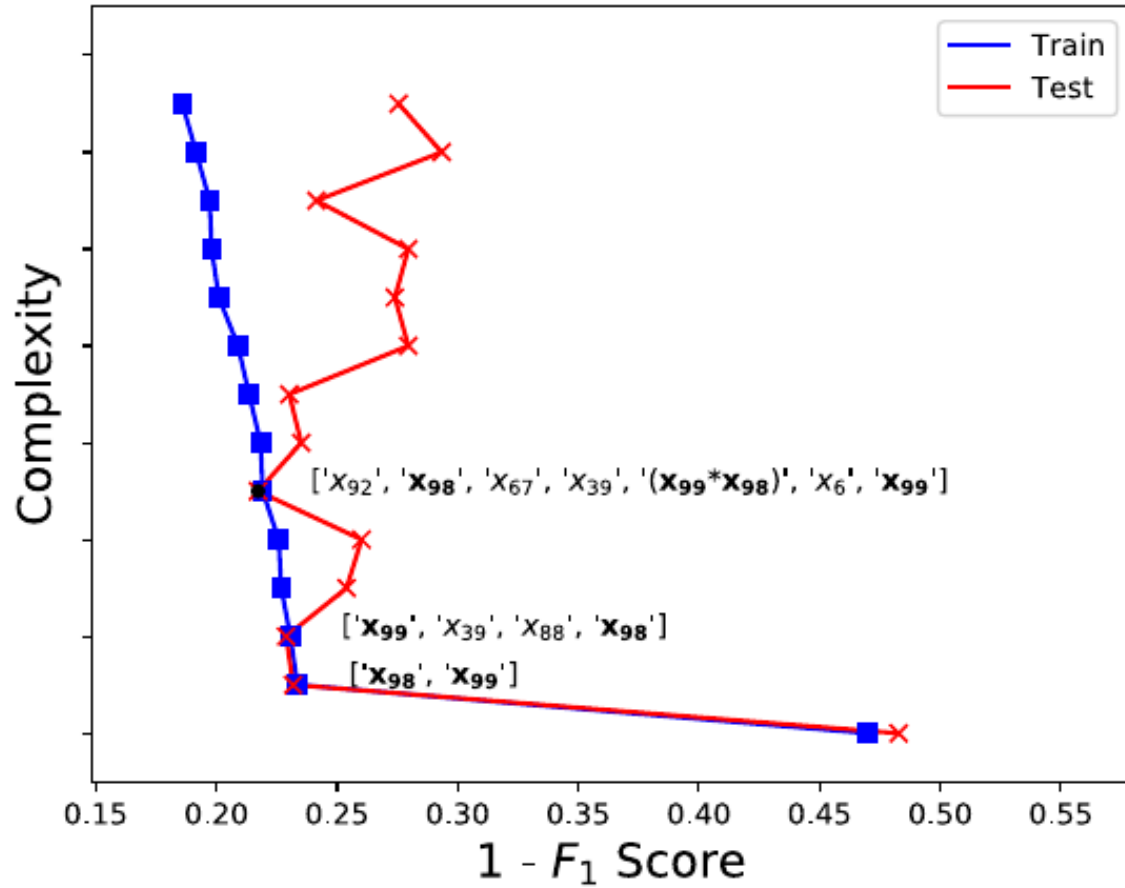
TABLE V  
COMPUTATION TIME (IN MINUTES)

Dataset	CSO	CCSO
Vehicle	<b>3.44</b>	4.41
WallRobot	92.12	<b>88.08</b>
German	<b>4.07</b>	4.66
GuesterPhase	375.01	<b>303.71</b>
Ionosphere	<b>1.21</b>	1.90
Chess	71.28	<b>37.68</b>
Movementlibras	<b>1.11</b>	1.42
Hillvalley	5.04	<b>4.78</b>
Musk1	1.81	<b>1.75</b>
USPS	240.33	<b>220.13</b>
Madelon	26.59	<b>25.35</b>
Isolet	<b>11.92</b>	15.80
MultipleFeatures	<b>18.80</b>	20.43
Gametes	14.93	<b>13.58</b>
QAR	<b>14.64</b>	22.51
QOT	308.07	<b>270.56</b>
COIL20	12.05	<b>11.20</b>
ORL	2.03	<b>1.94</b>
Bioresponse	<b>68.61</b>	77.73
RELATHE	<b>28.34</b>	29.86
GLIOMA	<b>0.51</b>	0.96
BASEHOCK	<b>55.81</b>	58.67
Gisette	<b>464.45</b>	487.06
Brain1	<b>0.93</b>	1.71





# Fronts



# Results

- Keep the results of the standard algorithm
- Do **keep all** the original results from **each run**.
- **Perform statistical significance tests**: T-test; Wilcoxon test; Friedman test
- **Do NOT delete** results unless they use too much memory, or they are wrong



# Codes and Programs

- **Backup** different versions of your codes using
  - Version control: Gitlab, bitbucket
- Make clear **documentation** of your codes
- Make clear **README** documentation
- **Organise your files:** a single directory for a project, containing all of the data, code, and results for your project
- **Following best programming practices:** choose file/method/variable name carefully, avoid hard coding, ...





**THANK**

**Questions?**