

Evaluation criteria for CEC 2024 competition and special session on numerical optimization considering accuracy and speed

Kangjia Qiao ^a, Xupeng Wen ^b, Xuanxuan Ban ^a, Peng Chen ^a,
Kenneth V. Price ^e, Ponnuthurai N. Suganthan ^d, Jing Liang ^{a,c},
Guohua Wu ^b, Caitong Yue ^a

a. School of Electrical and Information Engineering, Zhengzhou University,
Zhengzhou, 450001, China, qiaokangjia@gs.zzu.edu.cn, xxuanban@163.com,
ty1220899231@163.com

b. School of Traffic & Transportation Engineering, Central South University,
Changsha, 410073, China, wenxupeng@csu.edu.cn, guohuawu@csu.edu.cn

c. School of Electrical Engineering and Automation, Henan Institute of
Technology, Xinxiang, 453003, China, liangjing@zzu.edu.cn

d. KINDI Center for Computing Research, College of Engineering, Qatar
University, Doha, Box:2713, Qatar, p.n.suganthan@qu.edu.qa

e. Vacaville, CA, USA, pricekenneth459@gmail.com

Technical Report

November, 2023

CEC Reviewers' Requirement: CEC reviewers expect novel contributions in every submission. In addition, to be able to use the proposed U-score comparison approach, authors also need to include one or more algorithms taken from the literature. CEC paper is expected to include only the final results after exhausting the maximum number of FEs.

One submission is expected to address only one of the four cases. All results should be saved using high precision. Authors of accepted papers need to send the full results by email with a Readme.txt file. If you have any query, you can send an email to p.n.suganthan@qu.edu.qa and qiaokangjia@gs.zzu.edu.cn.

1. Introduction to the U-score approach

Traditionally, algorithm performance evaluation embraces one of the two complementary paradigms. The first, often referred to as the 'fixed target' scenario, records of the number of function evaluations (FEs) necessary for a trial to achieve a predetermined minimum function error value (Min_EV). The second, termed the 'fixed cost' scenario, records the function error value (EV) of a trial once it exhausts a stipulated maximum number of function evaluations (Max_FEs). Nonetheless, a notable lacuna has persisted in the assessment landscape, with a dearth of methodologies concurrently considering both Min_EV and Max_FEs.

The incorporation of both convergence accuracy and speed within the U-score approach [3] provides a comprehensive perspective on algorithmic performance, thereby facilitating effective comparative analyses and rankings across a multitude of algorithms by considering each run of each algorithms. Noteworthy is the fact that in the context of a binary competition involving just two algorithms, the U-score approach effectively simplifies to the Mann-Whitney U statistic.

Based on the U-score approach, we set up four groups of algorithmic ranking competitions, i.e., 1) Bound constrained single objective optimization problems. 2) Constrained single objective optimization problems. 3) Bound constrained multi-objective optimization problems. 4) Constrained multi-objective optimization problems.

2. Bound constrained single objective optimization problems (SOP)

U-score approach for single objective optimization problems [3]: To exemplify the U-score ranking method's confluence of convergence speed and accuracy on SOPs, an illustrative example is depicted in Figure 1. This figure portrays three ranking algorithms, designated as A1 to A3, each assigned a distinct color and shape. For each algorithm, four distinct runs were executed, yielding a total of 12 trials. These trials are stratified into two categories: 1) those that successfully converged to the Min_EV, target error value, and 2) those that fail to attain this target within the stipulated Max_FEs, concluded upon reaching the limit of Max_FEs. (PS: As the Min_EV is undefined as yet, authors are asked to execute to the Max_FEs and save results using high precision.)

In adherence to the stipulated procedures and criteria, the specific rankings of algorithms A1 to A3 are delineated within Figure 2, which presents the tabulated U-score results. To be specific, the scoring of algorithms A1, A2, and A3 is determined through the summation of their respective rankings. Evidently, algorithm A1 emerges as the victor, amassing a total score of 24 based on the U-score approach, thereby securing the topmost rank. In juxtaposition, A2 secures the second rank, and A3 the third. This outcome underscores the supremacy of algorithm A1, attributed to its swifter convergence velocity and its propensity to attain lower Min_EV.

Table 1

Results saved in “PaperID_FJ_Min_EV.mat” where J=1,2,3,...29 problems.

	Run 1	Run 2	Run 3	...	Run 25
Min_EV at Initialisation FEs					
Min_EV at 10*D FEs					
Min_EV at 20*D FEs					
...					
...					
Min_EV at Max_FEs					

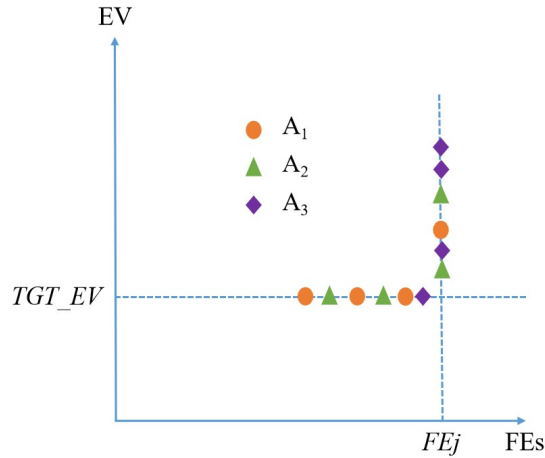


Figure 1: Three algorithms, A1–A3, run four trials each on an SOP. A single run terminates when it reaches Max_FEs. All trials’ results can be ordered from the best to the worst. TGT_EV and FE_j will be determined later.

Trial	●	▲	●	▲	●	◆	▲	◆	●	▲	◆	◆	SR	U-score
Ranks	12	11	10	9	8	7	6	5	4	3	2	1	78	
A1	12		10		8				4				34	24
A2		11		9			6			3			29	19
A3						7		5			2	1	15	5

The “correction factor” (*cf*) is $n(n + 1)/2 = 4 * 5/2 = 10$, where n denotes the number of trials. SR denotes the sum of ranks. The scores of algorithms are calculated by the “SR” minus the “*cf*” according to the U-score algorithm.

Figure 2: U-score ranks for algorithms A1, A2 and A3.

Test Problems: The 29 real-parameter numerical optimization problems with 30D in CEC2017 [1] are adopted as test problems. The codes can be downloaded from the website: <https://github.com/P-N-Suganthan/CEC2017-BoundContrained>.

Number of Trials/Problem: 25 independent runs. (Do not run many 25 runs to pick the best run).

Maximum Number of Function Evaluations: Max_FEs = 10000*D, where D is the dimensionality of the optimization problems.

Search Range: $[-100, 100]^D$

Population Size: You are free to have an appropriate population size to suit your algorithm while not exceeding the Max_FEs.

Sampling Points: The best EV (Error Value) around every $10 \cdot D$ evaluations will be recorded for each run. For example, the maximum number of function evaluations Max_FEs is $10000 \cdot D$, then 1000 EVs should be saved.

Target Error Values: The target error value, TGT_EV for each problem, will be determined after the competition. Hence, all algorithms should be executed until the Maximum number of Function Evaluations (Max_FEs) are consumed.

Algorithm Complexity: The evaluation of algorithm complexity requires the calculation of two indicators T_1 and T_2 , which are calculated as follows:

$$1) T_1 = (\sum_{i=1}^{29} t_i^1) / 29, t_i^1 \text{ is the computing time of 10000 evaluations for problem } i.$$

$$2) T_2 = (\sum_{i=1}^{29} t_i^2) / 29, t_i^2 \text{ is the complete computing time for the algorithm with 10000 evaluations for problem } i.$$

The complexity of the algorithm is reflected by: T_1 , T_2 and $(T_2 - T_1) / T_1$

Presentation of Results: The results can be saved in the form of Table 1, where Min_EV is the best error value of each run at each sampling point. The value should be recorded every $10 \cdot D$ FEs.

Thus, for each algorithm, 29 files should be zipped and sent to the organizers, where 29 represents the total number of test functions.

Note that all participants are allowed to improve their algorithms further after submitting the initial version of their papers until the final accepted paper submission deadline set by the conference. Authors are required to submit their results in the introduced format to the organizers after submitting the final version of paper as soon as possible.

3. Constrained single objective optimization problems

U-score approach for constrained single objective optimization problems (COPs): For constrained single objective optimization problems, the minimum objective function value f_{min} is the performance indicator at each sampling point, if the corresponding solution is feasible. The feasibility condition is defined in [5]. To exemplify the U-score ranking method's confluence of convergence speed and accuracy on COPs, the illustrative example is depicted in Figure 3. This figure portrays three ranking algorithms, designated as A1 to A3, each assigned a distinct color and shape. For each algorithm, four distinct runs were executed, yielding a total of 12 trials. These trials are stratified into two categories: 1) those that successfully converged to the TGT_f target value, and

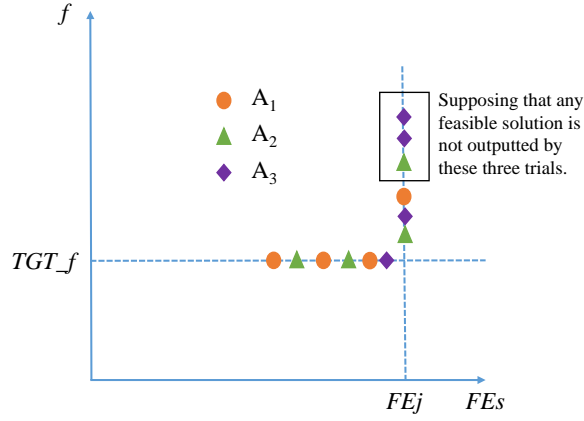


Figure 3: Three algorithms, A1–A3, run four trials each on a COP. A single run terminates when it reaches Max_FEs. TGT_f and FE_j will be determined later. All trial results can be ordered from the best to the worst.

Trial	●	▲	●	▲	●	◆	▲	◆	●	▲	◆	◆	SR	U-score
Ranks	12	11	10	9	8	7	6	5	4	3	2	1	78	
A1	12		10		8				4				34	24
A2		11		9			6			3			29	19
A3						7		5			2	1	15	5

The “correction factor” (cf) is $n(n + 1)/2 = 4 * 5/2 = 10$, where n denotes the number of trials. SR denotes the sum of ranks. The scores of algorithms are calculated by the “SR” minus the “ cf ” according to the U-score algorithm.

Figure 4: U-score ranks for CMOPs.

2) those that fail to attain this target within the stipulated Max_FEs. (PS: As the TGT_f is undefined as yet, authors are asked to execute to the Max_FEs and save results using high precision.)

In adherence to the stipulated procedures and criteria, the specific rankings of algorithms A1 to A3 are delineated within Figure 4, which presents the tabulated U-score results. To be specific, the scoring of algorithms A1, A2, and A3 is determined through the summation of their respective rankings. Evidently, algorithm A1 emerges as the victor, amassing a total score of 24 based on the U-score approach, thereby securing the topmost rank. In juxtaposition, A2 secures the second rank, and A3 the third. This outcome underscores the supremacy of algorithm A1, attributed to its swifter convergence velocity and its propensity to attain lower f value.

Please note that some algorithms might not output any feasible solution at a particular sampling point, that is, the entire population is infeasible. If two or more algorithms do not output any feasible solution at a sampling point, we will rank these trials based on the lowest value of the overall constraint violation (LCV) of the best solution. The algorithm with smaller LCV value is better. The formula for calculating LCV for a run of an algorithm is as follows:

$$LCV = \min: CV(P_i), i = 1, \dots, NP \quad (1)$$

where NP is the population size and $CV(P_i)$ is the overall constraint violation value of the i^{th} individual in the population P . Please note that when recording the LCV value in Table 2, CV should be directly calculated through the constraint functions, rather than the author's normalized CV result.

Test Problems: The 28 constrained real-parameter optimization problems with $30D$ in CEC2017 [5] are adopted as test problems. The code can be found from the PlatEMO (website: <https://github.com/BIM> or <https://github.com/P-N-Suganthan/CEC2017>).

Number of Trials/Problem: 25 independent runs.

Maximum Number of Function Evaluations: $\text{Max_FEs} = 20000 * D$, where D is the dimensionality of the optimization problems.

Population Size: You are free to have an appropriate population size to suit your algorithm while not exceeding the Max_FEs .

Sampling Points: The f_{min} values and LCV around every $10 * D$ evaluation will be recorded. For example, if the maximum number of function evaluations Max_FEs is $20000 * D$, then $2000 f_{min}$ values are recorded for trials with one or more feasible solutions. When the whole population is infeasible, the lowest LCV value of the population should be saved at the respective sampling points.

Target Error Values: The target error value will be determined after the competition. Hence, all algorithms should be executed until Maximum number of Function Evaluations (Max_FEs) are consumed.

Algorithm Complexity: The evaluation of algorithm complexity requires the calculation of two indicators T_1 and T_2 , which are calculated as follows:

- 1) $T_1 = (\sum_{i=1}^{28} t_i^1) / 28$, t_i^1 is the computing time of 10000 evaluations for problem i .
- 2) $T_2 = (\sum_{i=1}^{28} t_i^2) / 28$, t_i^2 is the complete computing time for the algorithm with 10000 evaluations for problem i .

The complexity of the algorithm is reflected by: T_1 , T_2 and $(T_2 - T_1) / T_1$

Presentation of Results: Save your results as shown in Table 2, in which the first entry is for the evaluation of the initial population. The cumulative FEs at each sampling point should be saved in the first column. Meanwhile, the corresponding f_{min} and LCV results should be saved in the second and third columns, respectively. So, for a function, one run requires one file in mat format. Please note that if no feasible solution exists at one sampling point, the f_{min} result should be expressed by "NaN".

Thus, for each algorithm, 28 files should be zipped and sent to the organizers, where 28 represents the total number of test functions.

Table 2

Results saved in "PaperID_CPJ.mat" where J=1,2,...,28 problems

FEs	Run1		Run2		...	Run25	
	f_{min}	LCV	f_{min}	LCV		f_{min}	LCV
at Initialisation FEs							
Sampling Point 1, FEs=1*10D							
Sampling Point 2, FEs=2*10D							
...							
Last Sampling Point, Max_FEs							

Note that all participants are allowed to improve their algorithms further after submitting the initial version of their papers until the final accepted paper submission deadline set by the conference. Authors are required to submit their results in the introduced format to the organizers after submitting the final version of paper as soon as possible.

Table 3
U-score ranks for MOEAs

Trial	●	●	●	●	●	●	●	●	●	●	●	●	●	●	●	●	SR	U-score ¹
Ranks	16	15	14	13	12	11	10	9	8	7	6	5	4	3	2	1	136	
A1	16		14						8	7							45	35
A2		15		13							6	5					39	29
A3					12			9					4		2		27	17
A4						11	10							3		1	25	15

¹ The "correction factor" cf is $n(n+1)/2 = 4 * 5/2 = 10$, where n denotes the number of trials. SR denotes the sum of ranks. The scores of algorithms are calculated by the " SR " minus the " cf " according to the U-score algorithm.

4. Bound constrained multi-objective optimization problems (MOPs)

U-score approach for unconstrained multi-objective optimization problems: For unconstrained multi-objective optimization problems, we use the Inverted Generational Distance (IGD) value as the performance indicator. To exemplify the U-score ranking method's confluence of convergence speed and accuracy on MOP, the illustrative example is depicted in Figure 5. This figure portrays four ranking algorithms, designated as A1 through A4, each assigned a distinct color. For each algorithm, four distinct runs were executed, yielding a total of 16 trials. These trials are stratified into two categories: 1) those that successfully converged to the TGT_IGD target value, and 2) those that, failing to attain this target within the stipulated Max_FEs, concluded upon reaching the limit of Max_FEs. (PS: As the TGT_IGD is undefined as yet, authors are asked to execute to the Max_FEs and save results using high precision.)

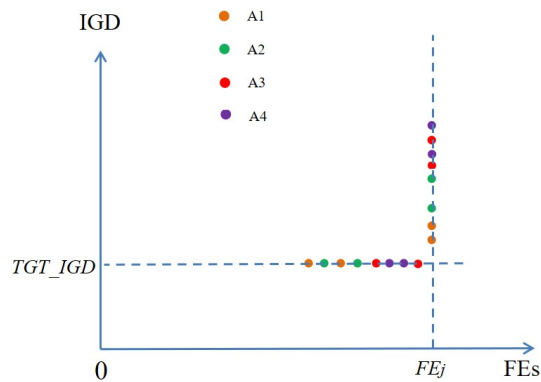


Figure 5: Four algorithms, A1–A4, run four trials each on an MOP. A run terminates when it reaches Max_FEs. TGT_IGD and FE_j will be determined later. All trial results can be ordered from best to worst.

Table 3 presents the tabulated U-score results, wherein the scoring of algorithms A1, A2, A3, and A4 is determined through the summation of their respective rankings. Evidently, algorithm A1 emerges as the victor, amassing a total score of 35 based on the U-score approach, thereby securing the topmost rank. In juxtaposition, A2 secures the second rank, A3 the third, and A4 the fourth. This outcome underscores the supremacy of algorithm A1, attributed to its swifter convergence velocity and its propensity to attain lower *IGD* values.

Test Problems: We adopt the benchmark of [2] including 10 multi-objective problems to rank the optimizers of MOPs without constraints.

Number of Trials/Problem: 30 independent runs.

Maximum Number of Function Evaluations: The maximum number of evaluations are set to 100000 for each function.

Pareto Front Size: The final PF (Front 1) is expected to have a size of 100. Compute *IGD* results using maximal 100 feasible individuals. The recommended population size is 100.

Sampling Points: The *IGD* values will be recorded once every 200 function evaluations. For example, if the maximum number of evaluations Max_FEs is 100000, then 500 *IGD* values are saved.

Target IGD Values: The target *IGD* value will be determined after the competition. Hence, all algorithms should be executed until Maximum number of Function Evaluations (Max_FEs) are consumed.

Encoding: If the algorithm requires encoding, then the encoding scheme should be independent of the specific problems and governed by generic factors such as the search ranges.

Algorithm Complexity: The evaluation of algorithm complexity requires the calculation of two indicators T_1 and T_2 , which are calculated as follows:

- 1) $T_1 = (\sum_{i=1}^{10} t_i^1)/10$, t_i^1 is the computing time of 10000 evaluations for problem i .
- 2) $T_2 = (\sum_{i=1}^{10} t_i^2)/10$, t_i^2 is the complete computing time for the algorithm with 10000 evaluations for problem i .

The complexity of the algorithm is reflected by: T_1 T_2 and $(T_2 - T_1)/T_1$

Presentation of Results: To compare and evaluate the algorithms participating in the competition, it is necessary that the authors email the results as shown in Table 4 to the organizers after submitting the final version of the accepted paper.

According to Table 4, 501 *IGD* values for each of the 30 runs are required for each problem. For example, the results of PaperID for problem RCMJ, the file name should be “PaperID RCMJ IGD.txt”, where Inverted Generational Distance values are saved, respectively. Thus, $10 * 30 =$

Table 4Results saved in "PaperID RCMJ IGD.txt" where $J=1,2,\dots,10$ problems

	Run 1	Run 2	Run 3	...	Run 30
<i>IGD</i> at Initialisation FEs					
<i>IGD</i> at Sampling Point 1					
<i>IGD</i> at Sampling Point 2					
...					
...					
<i>IGD</i> at Sampling Point 500, 100K FEs					

300 files should be zipped and sent to the organizers, where 10 represents the total number of test functions, and 30 represents the number of trials per problem.

Note that all participants are allowed to improve their algorithms further after submitting the initial version of their papers until the final accepted paper submission deadline set by the conference. Authors are required to submit their results in the introduced format to the organizers after submitting the final version of paper as soon as possible. In summary, the results should be saved as shown in Table 4.

5. Constrained multi-objective optimization problems (CMOPs)

U-score approach for constrained multi-objective optimization problems: For constrained multi-objective optimization problems, we introduce the Inverted Generational Distance (*IGD*) values as an indicator. To exemplify the U-score ranking method's confluence of convergence speed and accuracy on CMOPs, the illustrative example is depicted in Figure 6. This figure portrays three ranking algorithms, designated as A1 to A3, each assigned a distinct color and shape. For each algorithm, four distinct runs were executed, yielding a total of 12 trials. These trials are stratified into two categories: 1) those that successfully converged to the TGT_IGD target value, and 2) those that fail to attain this target within the stipulated Max_FEs. (PS: As the TGT_IGD is undefined as yet, authors are asked to execute to the Max_FEs and save results using high precision.)

In adherence to the stipulated procedures and criteria, the specific rankings of algorithms A1 to A3 are delineated within Figure 7, which presents the tabulated U-score results. To be specific, the scoring of algorithms A1, A2, and A3 is determined through the summation of their respective rankings. Evidently, algorithm A1 emerges as the victor, amassing a total score of 24 based on the U-score approach, thereby securing the topmost rank. In juxtaposition, A2 secures the second rank, and A3 the third. This outcome underscores the supremacy of algorithm A1, attributed to its swifter convergence velocity and its propensity to attain lower *IGD* values.

Please note that some algorithms might not output any feasible solution on one run, that is, the entire population is infeasible. If at least two trials do not output any feasible solutions as such in Figure 6, we will rank these trials based on the mean value of overall constraint violation value of the populations (*MCV*). The trial with a smaller *MCV* value is better. The formulate of calculating *MCV* is as follows:

$$MCV = \frac{\sum_{i=1}^{PF} CV(P_i)}{PF} \quad (2)$$

where *PF* is the required final front-1 size and $CV(P_i)$ is the constraint violation value of the i^{th} individual in the population *P*. Please note that when recording the *MCV* value in Table 5, *CV* should be **directly calculated through the constraint functions**, rather than the author's normalized *CV* result.

Test Problems: The latest constrained multiobjective optimization problems with scalable decision space constraints (SDC problems) [4] are adopted as test problems. SDC benchmark contains 15 problems. The codes can be downloaded from the website¹.

Number of Trials: 30 independent runs.

Maximum Number of Function Evaluations: The maximum number of evaluations are set to 200000 for each function.

¹<https://github.com/cilabzzu/Codes/blob/main/SDC>, which should be ran on PlatEMO (website: <https://github.com/BIMK/PlatEMO>). Please note that the latest PlatEMO has contained the codes of SDC functions, and you can directly use these codes on PlatEMO.

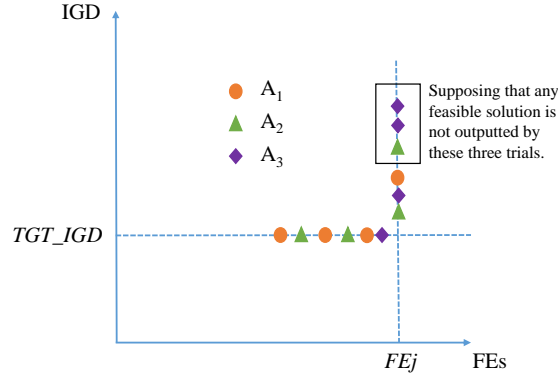


Figure 6: Three algorithms, A1–A3, run four trials each on a CMOP. A single run terminates when reaches Max_FEs. TGT_IGD and FE_j will be determined later. All trial results can be ordered from the best to the worst.

Trial	●	▲	●	▲	●	◆	▲	◆	●	▲	◆	◆	SR	U-score
Ranks	12	11	10	9	8	7	6	5	4	3	2	1	78	
A1	12		10		8				4				34	24
A2		11		9			6			3			29	19
A3						7		5			2	1	15	5

The “correction factor” (cf) is $n(n + 1)/2 = 4 * 5/2 = 10$, where n denotes the number of trials. SR denotes the sum of ranks. The scores of algorithms are calculated by the “SR” minus the “ cf ” according to the U-score algorithm.

Figure 7: U-score ranks for CMOPs.

Pareto Front Size: The final PF (i.e. Front 1) is expected to have a size of 100. Compute IGD results using maximal 100 feasible individuals. The recommended population size is 100.

Parameter Setting: The dimension is set to 30 for each SDC function.

Sampling Points: The IGD values will be recorded once every 200 function evaluations. For example, if the maximum number of evaluations Max_FEs is 200000, then 1000 IGD values are saved.

Target IGD Values: The target IGD value will be determined after the competition. Hence, all algorithms should be executed until the Maximum number of Function Evaluations (Max_FEs) are consumed. Please note that the minimal IGD value is unknown for multiobjective optimization problems. So, the mean or median IGD value of all trials from all algorithms participating in the competition will be set as the target IGD value.

Algorithm Complexity: The evaluation of algorithm complexity requires the calculation of two indicators T_1 and T_2 , which are calculated as follows:

$$1) T_1 = (\sum_{i=1}^{15} t_i^1)/15, t_i^1 \text{ is the computing time of 10000 evaluations for problem } i.$$

Table 5

Results saved in "PaperID_SDCJ.mat" where J=1,2,...,15 problems.

	Run1		Run2		...	Run 30	
	<i>IGD</i>	<i>MCV</i>	<i>IGD</i>	<i>MCV</i>		<i>IGD</i>	<i>MCV</i>
at initialization FEs							
Sampling point 1							
Sampling point 2							
...							
...							
Sampling point 1000							

2) $T_2 = (\sum_{i=1}^{15} t_i^1)/15$, t_i^2 is the complete computing time for the algorithm with 10000 evaluations for problem i .

The complexity of the algorithm is reflected by: T_1 , T_2 and $(T_2 - T_1)/T_1$.

Presentation of Results: To compare and evaluate the algorithms participating in the competition, it is necessary that the authors email the results in the format as shown in Table 5 to the organizers, after submitting the final version of the accepted papers.

In Table 5, at each sampling point, i.e. every 200 FEs, *IGD* and *MCV* should be computed and saved in the second and third columns, respectively. Please note that if no feasible solution exists at a sampling point, the *IGD* result should be expressed by "NaN", while *MCV* value should be recorded.

For one algorithm, 15 files in .mat format (one for each problem) should be zipped using the paper number as the file name and sent to the organizers.

Note that all participants are allowed to improve their algorithms further after submitting the initial version of their papers until the final accepted paper submission deadline set by the conference. Authors are required to submit their final results in the prescribed format to the organizers after submitting the final version of paper as soon as possible.

References

- [1] Awad, N., Ali, M., Liang, J., Qu, B., Suganthan, P., 2016. Problem definitions and evaluation criteria for the cec 2017 special session and competition on single objective bound constrained real-parameter numerical optimization, in: Technical Report. Nanyang Technological University Singapore, pp. 1–34.
- [2] Li, H., Deb, K., Zhang, Q., Suganthan, P.N., Chen, L., 2019. Comparison between moea/d and nsga-iii on a set of novel many and multi-objective benchmark problems with challenging difficulties. *Swarm and Evolutionary Computation* 46, 104–117.
- [3] Price, K.V., Kumar, A., Suganthan, P.N., 2023. Trial-based dominance for comparing both the speed and accuracy of stochastic optimizers with standard non-parametric tests. *Swarm and Evolutionary Computation*, 101287, April .
- [4] Qiao, K., Liang, J., Yu, K., Yue, C., Lin, H., Zhang, D., Qu, B., 2023. Evolutionary constrained multiobjective optimization: Scalable high-dimensional constraint benchmarks and algorithm. *IEEE Transactions on Evolutionary Computation* , 1–1doi:10.1109/TEVC.2023.3281666.
- [5] Wu, G., Mallipeddi, R., Suganthan, P.N., 2017. Problem definitions and evaluation criteria for the cec 2017 competition on constrained real-parameter optimization. National University of Defense Technology, Changsha, Hunan, PR China and Kyungpook National University, Daegu, South Korea and Nanyang Technological University, Singapore, Technical Report .