

Finding a Good Length for Automatically Generated Tests

TABLE I
PYTHON SUBJECTS

SUT	LOC	Source
AVL	225	TSTL example [6]
bidict	569	http://pythonhosted.org/bidict/home.html
CParser	5033	https://github.com/albertz/PyCParser
pyfakefs	2642	https://github.com/jmcgeheeiv/pyfakefs
redis-py	2722	https://github.com/andymccurdy/redis-py
python-rsa	1597	https://github.com/sybrenstuvcl/python-rsa
simplejson	2811	https://simplejson.readthedocs.io/en/latest/
sortedcontainers	2017	http://www.grantjenks.com/docs/sortedcontainers/
SymPy	227959	http://www.sympy.org/en/index.html

I. INTRODUCTION

There are a lot of papers out there on this topic [1]–[5].

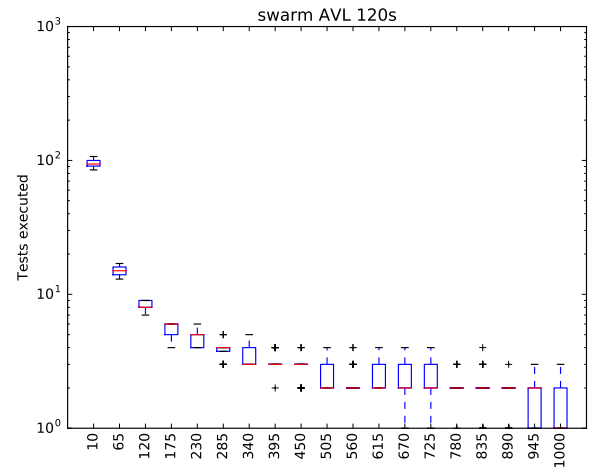
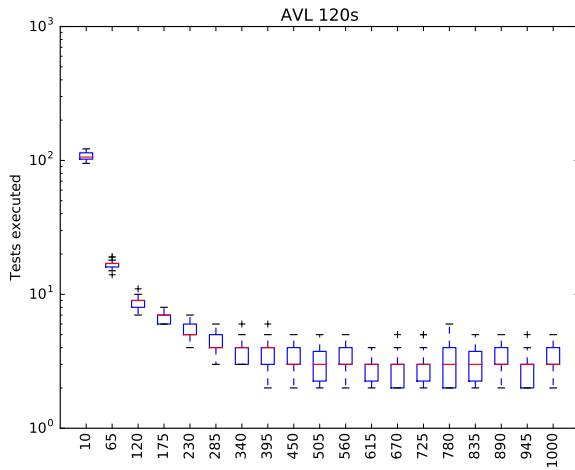
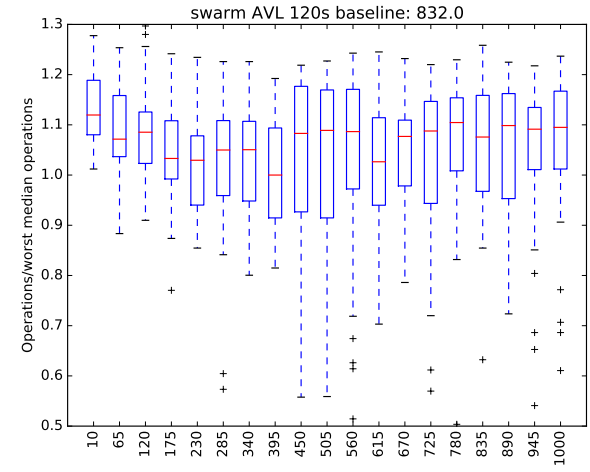
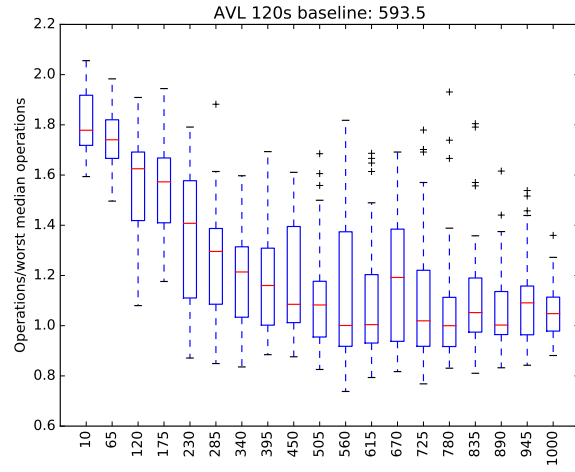
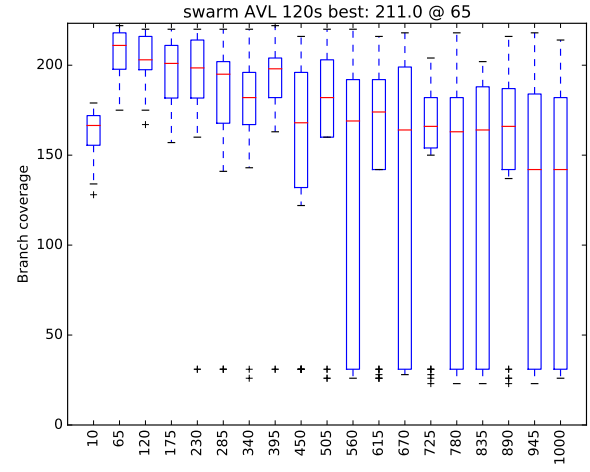
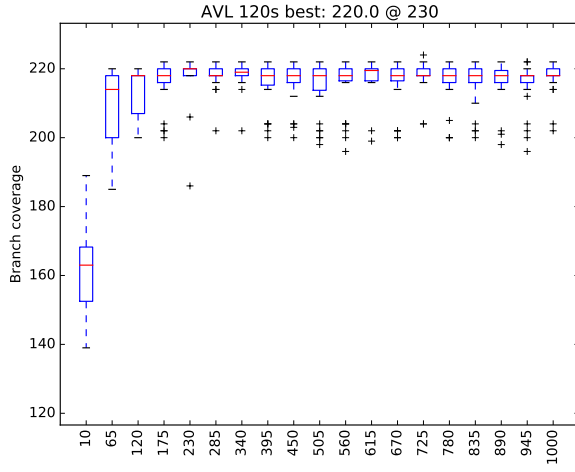
II. SOME PYTHON RESULTS

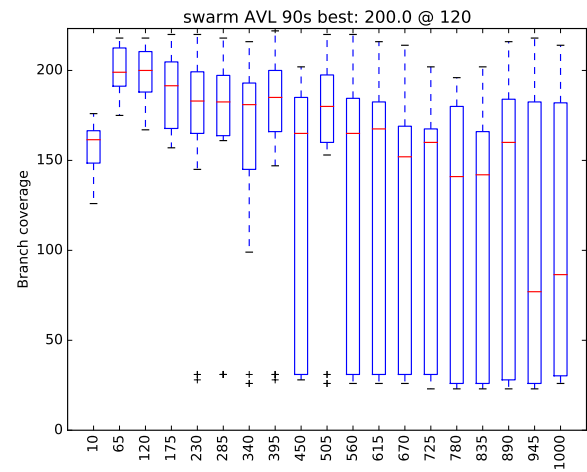
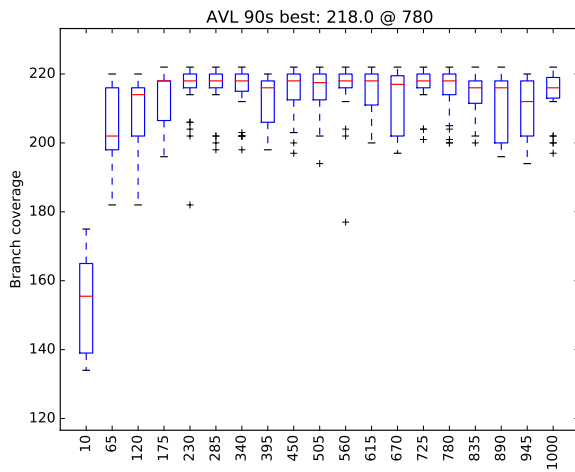
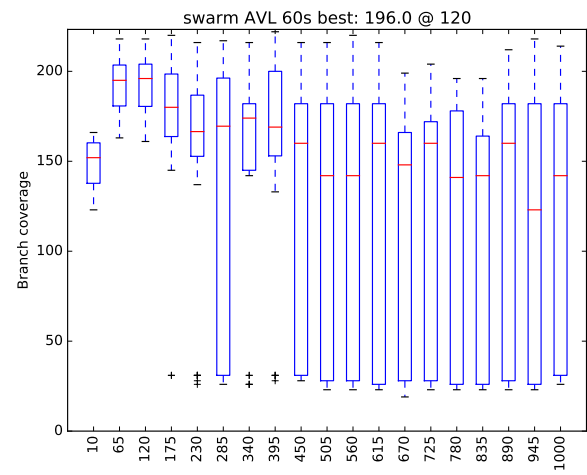
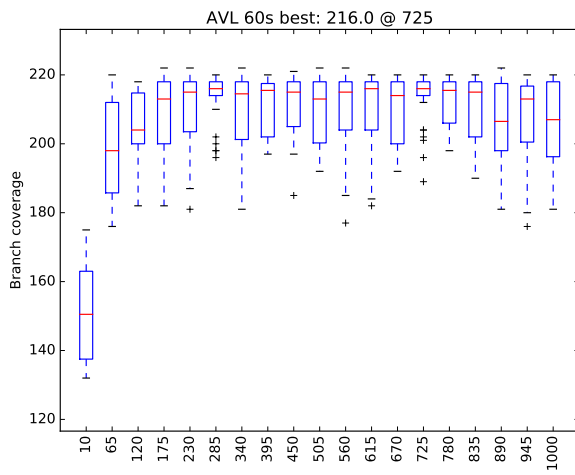
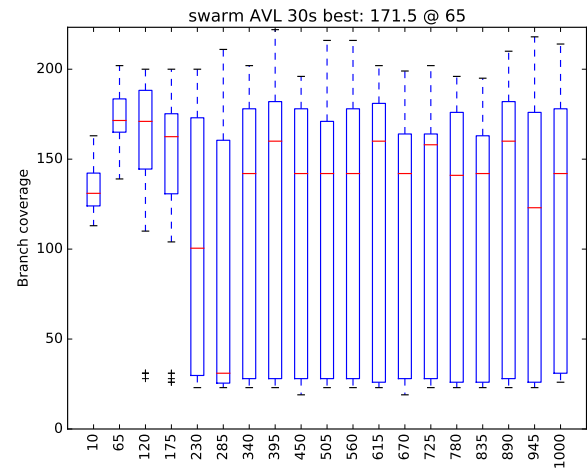
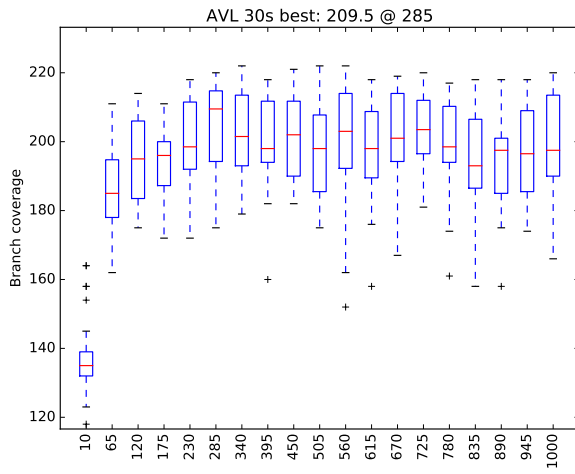
Table I describes the Python subjects, tested using the TSTL [7], [8] harnesses included in the TSTL distribution. Of the eight subjects, one is a toy container class with an injected realistic fault, and the other seven are real-world Python libraries, with a large number of GitHub stars or high ranking in pip. Results are shown for both standard random testing and swarm testing [9], using 120 seconds of test time. While the change to swarm testing sometimes changes the effectiveness curve, shapes at 30 seconds and 60 seconds are very similar, despite changes in overall coverage obtained.

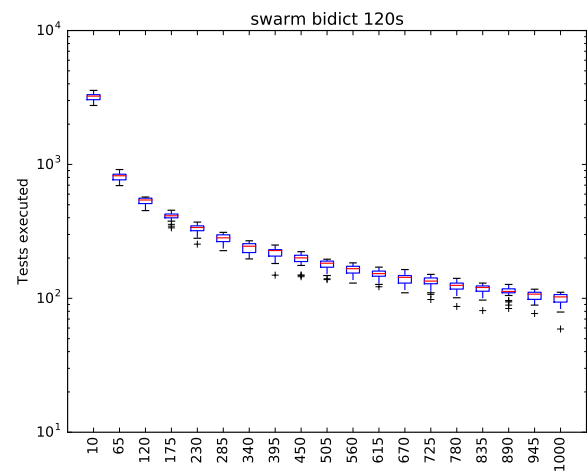
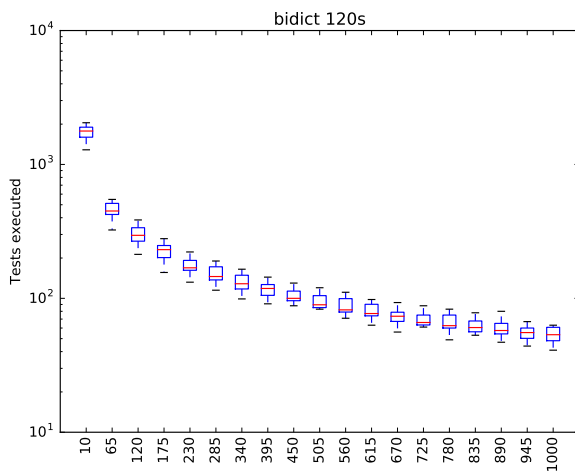
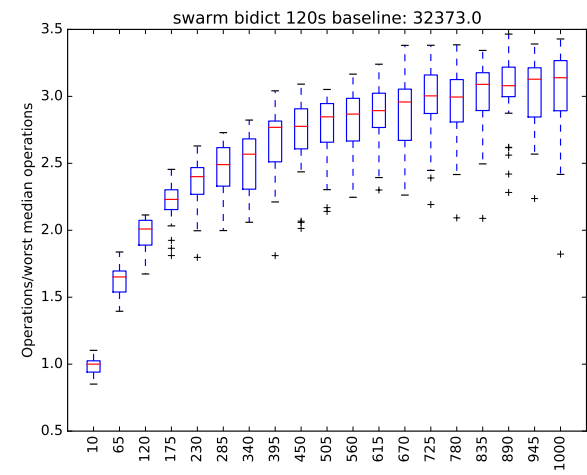
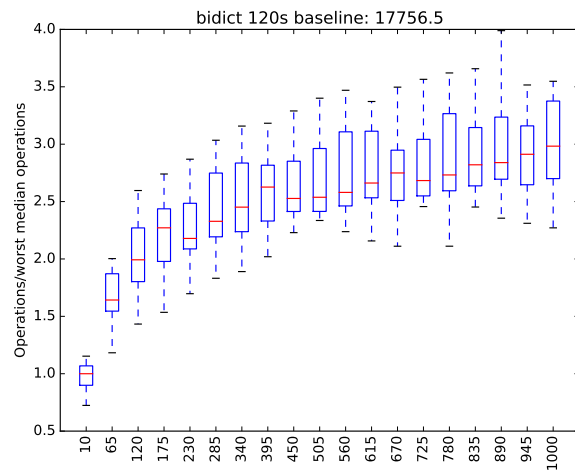
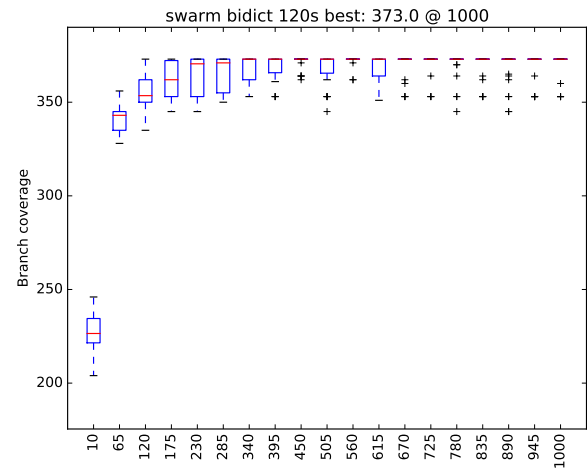
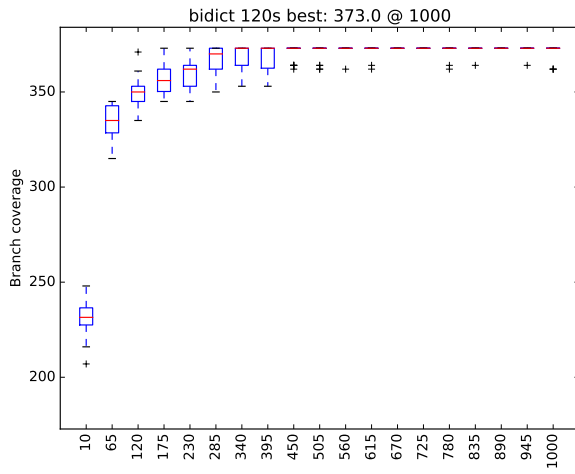
III. RELATED WORK

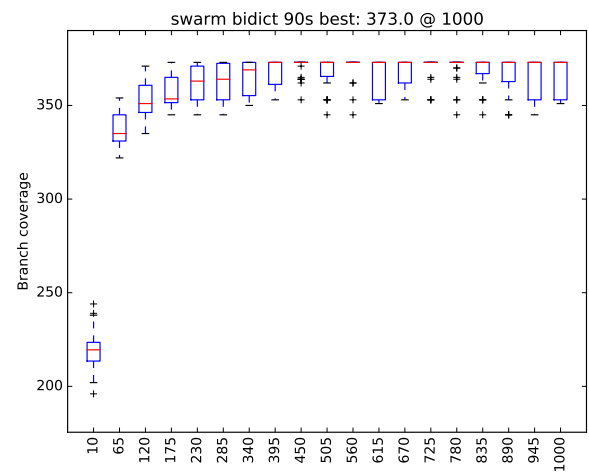
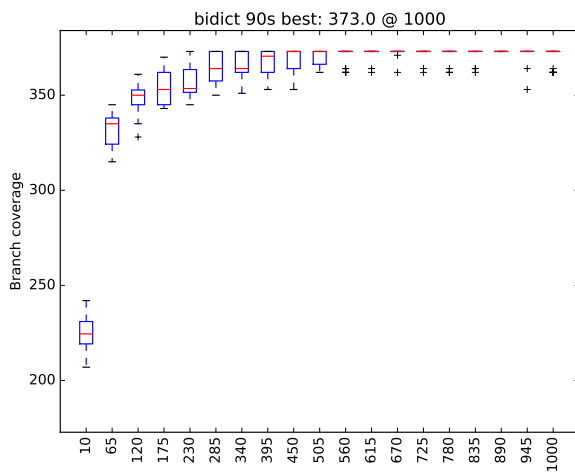
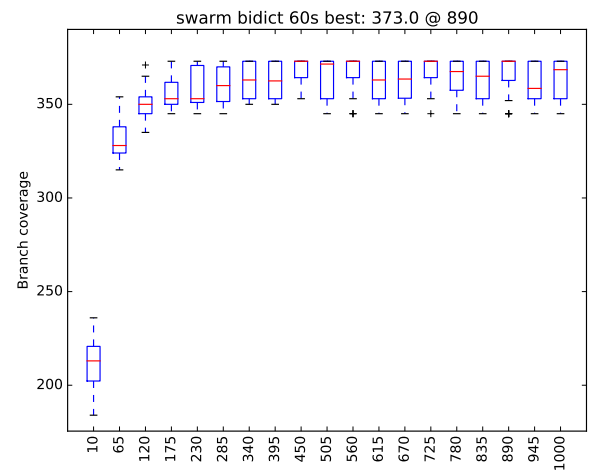
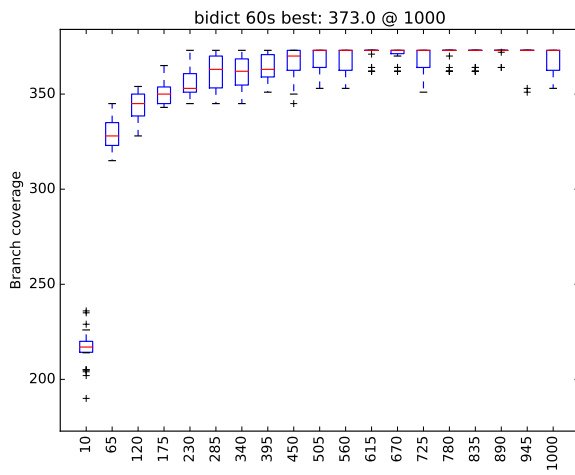
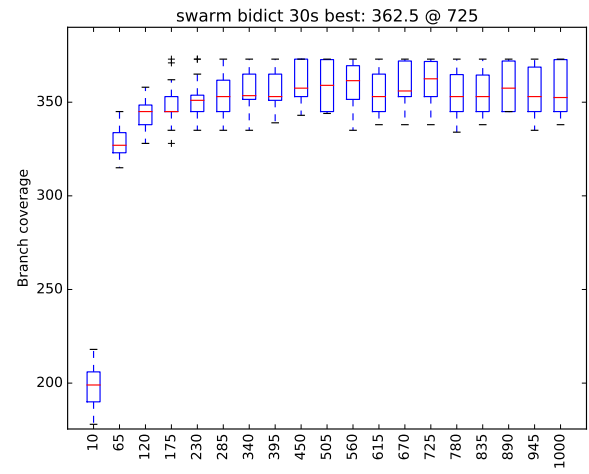
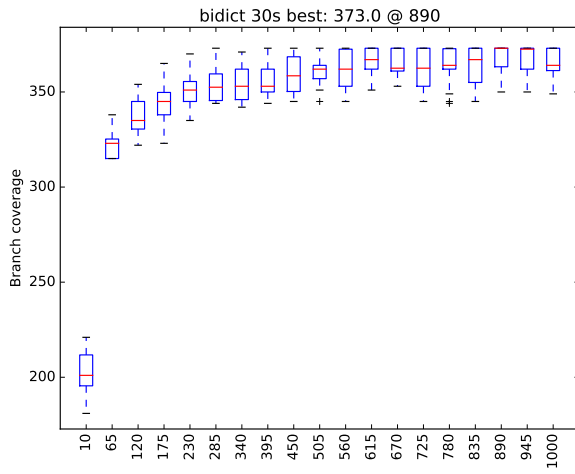
REFERENCES

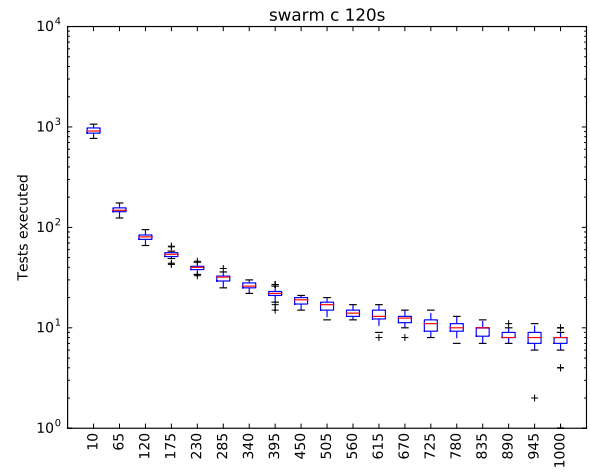
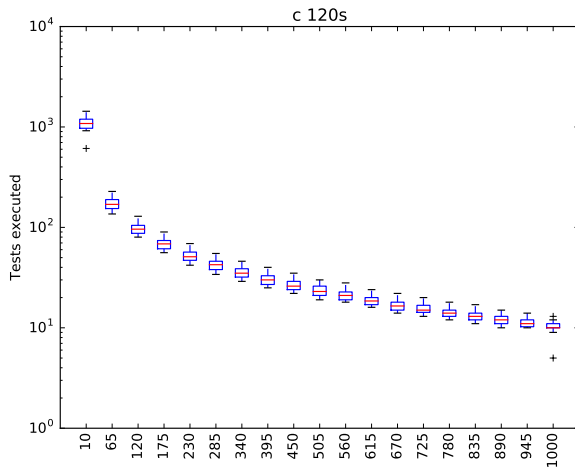
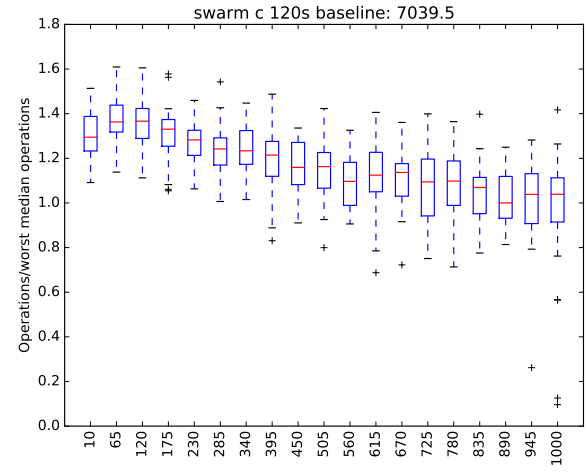
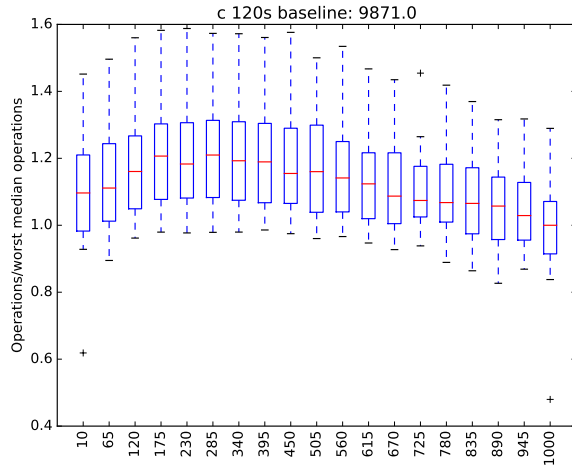
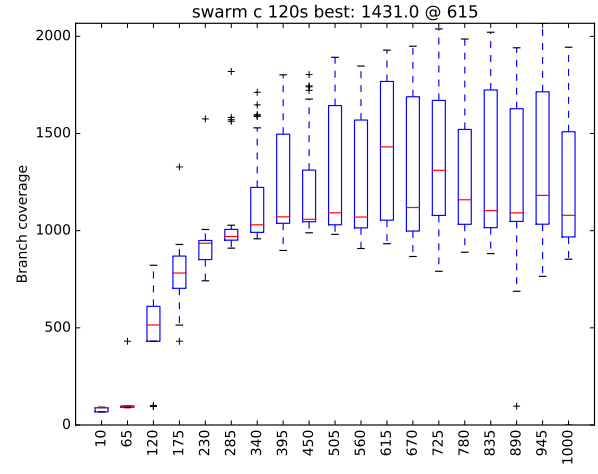
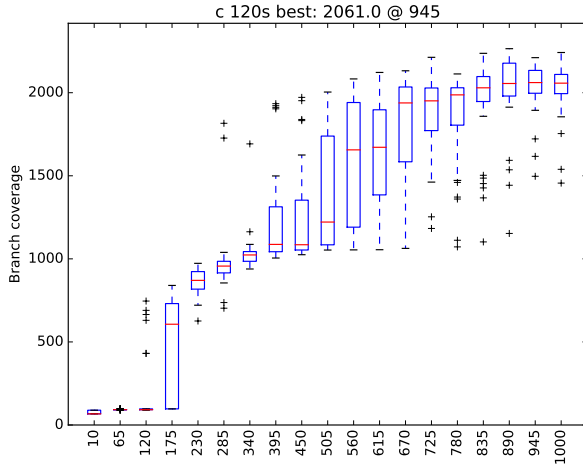
- [1] J. H. Andrews, A. Groce, M. Weston, and R.-G. Xu, “Random test run length and effectiveness,” in *Automated Software Engineering*, 2008, pp. 19–28.
- [2] A. Arcuri, “Longer is better: On the role of test sequence length in software testing,” in *Third International Conference on Software Testing, Verification and Validation, ICST 2010, Paris, France, April 7-9, 2010*. IEEE Computer Society, 2010, pp. 469–478. [Online]. Available: <https://doi.org/10.1109/ICST.2010.16>
- [3] G. Fraser and A. Arcuri, “It is not the length that matters, it is how you control it,” in *Fourth IEEE International Conference on Software Testing, Verification and Validation, ICST 2011, Berlin, Germany, March 21-25, 2011*. IEEE Computer Society, 2011, pp. 150–159. [Online]. Available: <https://doi.org/10.1109/ICST.2011.54>
- [4] A. Arcuri, “A theoretical and empirical analysis of the role of test sequence length in software testing for structural coverage,” *IEEE Trans. Software Eng.*, vol. 38, no. 3, pp. 497–519, 2012. [Online]. Available: <https://doi.org/10.1109/TSE.2011.44>
- [5] G. Fraser and A. Arcuri, “Handling test length bloat,” *Softw. Test., Verif. Reliab.*, vol. 23, no. 7, pp. 553–582, 2013. [Online]. Available: <https://doi.org/10.1002/stvr.1495>
- [6] user1689822, “python AVL tree insertion,” <http://stackoverflow.com/questions/12537986/python-avl-tree-insertion>.
- [7] J. Holmes, A. Groce, J. Pinto, P. Mittal, P. Azimi, K. Kellar, and J. O’Brien, “TSTL: the template scripting testing language,” *International Journal on Software Tools for Technology Transfer*, 2017, accepted for publication.
- [8] A. Groce and J. Pinto, “A little language for testing,” in *NASA Formal Methods Symposium*, 2015, pp. 204–218.
- [9] A. Groce, C. Zhang, E. Eide, Y. Chen, and J. Regehr, “Swarm testing,” in *International Symposium on Software Testing and Analysis*, 2012, pp. 78–88.



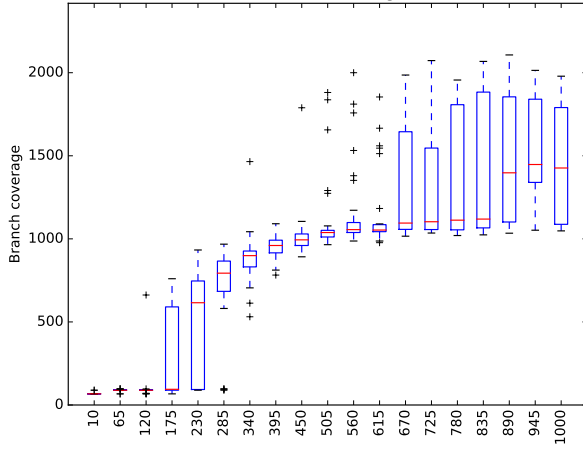




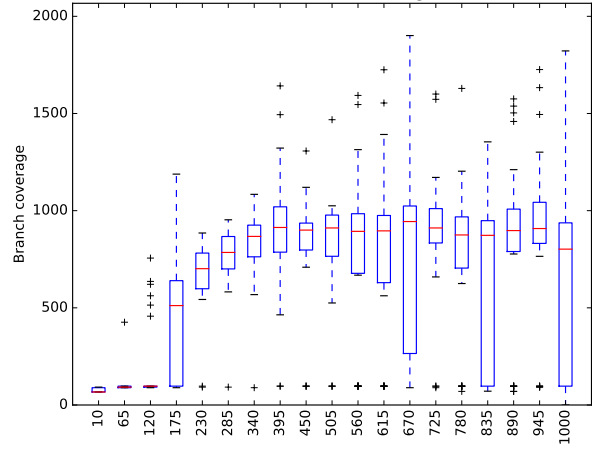




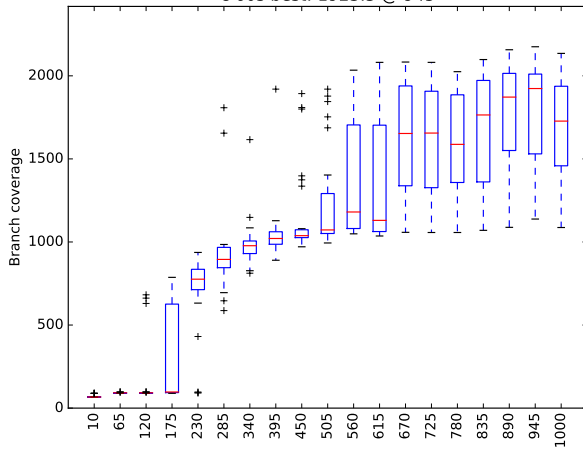
c 30s best: 1447.5 @ 945



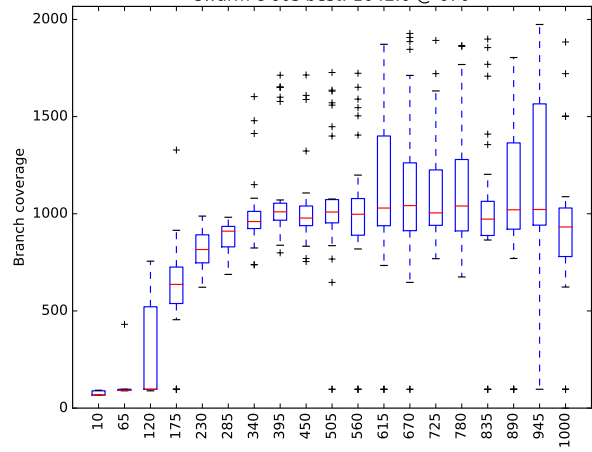
swarm c 30s best: 944.0 @ 670



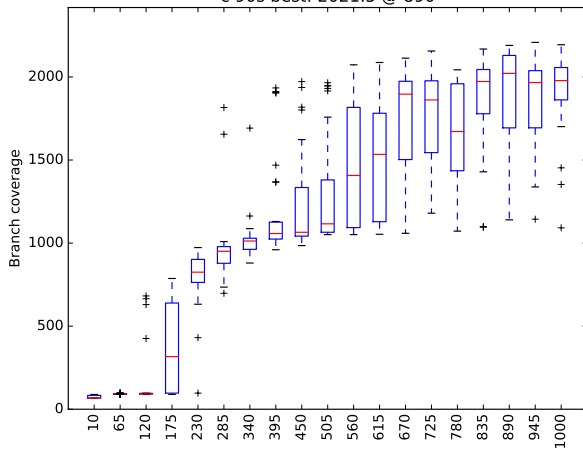
c 60s best: 1923.5 @ 945



swarm c 60s best: 1042.0 @ 670



c 90s best: 2021.5 @ 890



swarm c 90s best: 1324.0 @ 615

