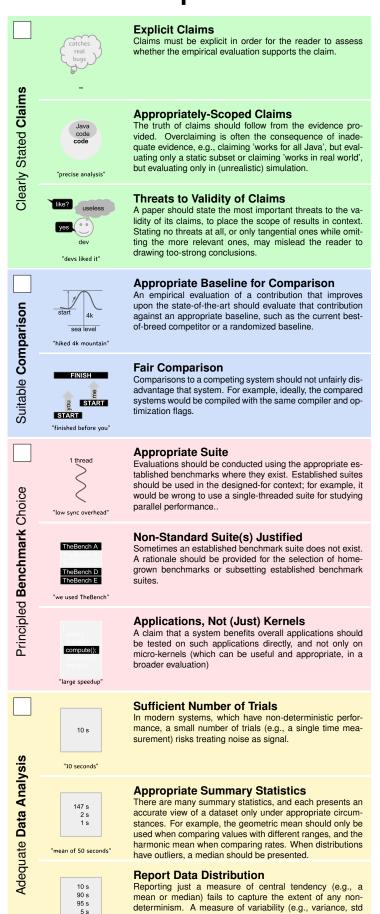
# **Empirical Evaluation Checklist** (alpha version)



deviation, quantiles) and/or confidence intervals help to un-

derstand the distribution of the data.

"50 seconds'



## **Direct or Appropriate Proxy Metric**

If the most relevant evaluation metric is not (or cannot be) measured directly, the proxy metric used instead must be well justified. For example, a reduction in cache misses is not an appropriate proxy for actual end-to-end performance or energy consumption.

## **Measures All Important Effects**

The costs and benefits of a technique may be multi-faceted. All facets should be considered, both costs and benefits. For example, compiler optimizations may speed up programs but at the cost of drastically increasing compile



"sped up apache"

## **Sufficient Information to Repeat**

Experiments should be described in sufficient detail to be repeatable. All parameters (including default values) should be included, as well as all version numbers of software, and full details of hardware platforms.

# Appropriate and Clear Experimental Design

**Reasonable Platform** 

The evaluation should be on a platform that can reasonably be said to match the claims. For example, a claim that relates to performance on mobile platforms should not have an evaluation performed exclusively on server.



**Explores Key Design Parameters** 

Key parameters should be explored over a range to evaluate sensitivity to their settings. Examples include the size of the heap when evaluating garbage collection and the size of caches when evaluating a locality optimization.



## Open Loop in Workload Generator

Load generators for transaction-oriented systems should not be gated by the rate at which the system responds. Rather, the load generator should be 'open loop', generating work independent of the performance of the system



19+3 22 7+1 8

test 19+3 22√ 7+1 8√

"perfect"

### **Cross-Validation Where Needed**

When a system aims to be general but was developed by training on or close consideration of specific examples, it is essential that the evaluation explicitly perform crossvalidation, so that the system is evaluated on data distinct from the training set.



## **Comprehensive Summary Results**

Appropriate statistics should be used to characterize the full range of results, not just the most favorable values, which may be outliers. For example, it is not appropriate to summarize speedups of 4%, 6%, 7%, and 49% as 'up to 49%'.

"have up to 4 leaves"



"B is much faster'

**Axes Include Zero** 

A truncated graph (with an axis not including zero) can exaggerate the importance of a difference. While 'zooming' in to the interesting range of an axis can potentially aid exposition, there is a significant risk that this is misleading (especially if it is not immediately clear that the axis is trun-



sped up B lots'

Appropriate Presentation of Results

**Ratios Plotted Correctly** 

When ratios such as speedups and slowdowns are plotted, the size of the bars must be linearly/logarithmically proportional to the change. When shown on the same linear scale, results are visually distorted by 1/r, where r is the ratio. This misleading effect can be avoided either by using a log scale or by normalizing to the lowest (highest) value.



## **Appropriate Level of Precision**

The number of significant digits should reflect the precision of the experiment. Reporting improvements of '49.9%' when the experimental error is +/- 1% is an example of misstated precision, misleading the reviewer's understanding of the significance of the rest.