

# Evaluating Security Products with Clinical Trials\*

Anil Somayaji    Yiru Li  
Carleton Computer Security Lab  
Ottawa, Ontario, Canada  
*{soma,yiruli}@ccsl.carleton.ca*

Hajime Inoue  
ATC-NY  
Ithaca, NY USA  
*hinoue@atc-nycorp.com*

José M. Fernandez  
École Polytechnique Montréal  
Montréal, Quebec, Canada  
*jose.fernandez@polymtl.ca*

Richard Ford  
Florida Insitute of Technology  
Melbourne, Florida, USA  
*rford@fit.edu*

June 2, 2009

## Abstract

One of the largest challenges faced by purchasers of security products is evaluating their relative merits. While purchasers can get reliable information on characteristics such as run-time overhead, user interface, and support quality, the actual level of protection provided by different security products is mostly unranked—or, worse yet, ranked using criteria that do generally reflect their performance in practice. Even though researchers have been working on improving testing methodologies, given the complex interactions of users, uses, evolving threats, and different deployment environments, there are fundamental limitations on the ability of lab-based measurements to determine real world performance. To address these issues, we propose an alternative evaluation method, computer security clinical trials. In this method, security products are deployed in randomly selected subsets of targeted populations and are monitored to determine their performance in practice. We believe that clinical trials can provide solid evidence of the efficacy of security products, much as they have in the field of medicine.

## 1 Introduction

The Internet is a dangerous place for users. As the reach of the network has increased, it has brought with it not only access to vast collections of data but also fraud and compromise. According to several reports [1], users are at more risk of attack than ever before. Furthermore, attackers are increasingly sophisticated, adapting quickly to new technologies and countermeasures, and nimbly morphing strategies to maximize payoffs. While the security industry has mounted a

---

\*Technical Report TR-09-05, School of Computer Science, Carleton University

valiant effort, we face a situation where our best efforts seem to be inadequate: attacks continue to occur, and users continue to be compromised.

Perhaps the scariest part of this situation is that we don't completely understand why we are failing. We have identifiable problems: unapplied patches, out-of-date malware signatures, poorly written software, complacent users. . . security experts can pontificate at length regarding the weaknesses of current systems. However, moving from this subjective, qualitative list to more concrete evaluations is difficult. Is patching *more* important than updating malware signatures? If so, how risky are delayed updates? And, more importantly, what defenses work in the field, and which ones do not? It is relatively easy to decide whether a defense could stop an attack; it is quite another to say that it will stop that attack in practice—particularly when attackers are given time to adapt and users are given the opportunity to invalidate the defense.

The key reason we don't know what defensive techniques work is that we can't reliably evaluate their relative merits *it situ*. Virus scanners perform very similarly in most lab tests—often the “best” solutions differ by fractions of one percent in overall results. Firewalls are compared and sold based upon features and speed, not security. And security experts regularly give advice such as “use strong passwords” and “turn off JavaScript” that place a huge burden on users yet, in many common situations, are completely inadequate against the actual threat. If security experts don't know what are the best security products, and we don't know what the best security advice to give to non-experts, is it any surprise that we have poorly secured systems?

While lab-based evaluations are important, we believe we must move beyond them if we are to make significant strides in improving the security of the Internet. Specifically, we must learn what works best *on deployed systems*. Note that “what works” is not the same as “what could work.” For example, usability studies can identify problems that could arise in deployment, such as difficulties in firewall configuration or confusion over messages from an antivirus scanner. Ultimately, though, we don't care about usability as determined in the laboratory—we care about actual use: Do administrators misconfigure firewalls in practice? How often does user confusion over proper virus scanner use actually lead to compromise?

To measure the use of security technologies in real-world circumstances, we have to account for how a given technology will interact with a huge variety of software, systems, users, uses, and attack profiles. The full complexity of the computational world cannot be captured in any lab setting or theoretical model—there are too many variables, and many of them change over timeframes (months or years) that cannot be practically measured in a laboratory setting using humans. As an alternative, we propose that the performance of security technologies be measured “in the field.” Specifically, we propose that security technologies be tested using the same methodology as used in medical *clinical trials*. In essence, we propose that we use the same measures of outcome, side effects, and user tolerance and compliance that regulatory bodies use to demonstrate that the benefit of a drug or medical device outweighs its risks. Clinical trials come in many forms depending upon the specific questions they are designed to address; what they all have in common, though, is that the test subjects live in the “real” world, not a laboratory.

Clinical trials were originally developed because medical practitioners faced challenges analogous to those faced by today's security professionals: they knew a lot about health problems, but they didn't know what worked to prevent or fix them. Clinical trials provided a methodology for separating “snake oil” from penicillin. As we will explain, clinical trials have a number of limitations as a testing methodology; our hope, though, is that clinical trials of security technologies will allow us to separate ineffective and dangerous technologies from those that provide significant

security benefits to real individuals and organizations.

The rest of this paper proceeds as follows. First, Section 2 discusses the evaluation problem in computer security. We briefly describe medical clinical trials in Section 3. In Section 4, we present our proposal for computer security clinical trials and outline a trial could be designed to evaluate antivirus products. Section 5 discusses a few significant potential objections to our proposal. Section 6 concludes.

## 2 Computer Security Problems

The evaluation problem exists broadly in computer security, for both academic research and commercial products. The most egregious type of improperly evaluated security technology is often referred to as “snake oil” [5]. The ultimate question in computer security evaluation is, how do we differentiate effective security mechanisms from such quackery, particularly in the eyes of a lay audience?

This problem is surprisingly hard. The key issue is that the current best-of-breed commercial systems are almost always unable to detect the most recent threats. This limitation arises from the current threat landscape. New threats emerge much more frequently than before, and meanwhile some of them, such as the Storm Worm [4], aim for economic profits and use very complex technologies in order to bypass security mechanisms. Even though many security companies have started using more flexible techniques such as heuristics to respond to new threats, in this arms race attackers always have an important advantage over security people—the public availability of security products. Highly-skilled attackers can keep modifying their newly created malicious codes until they can bypass all current defenses. The ease of achieving this has been proved by race-to-zero [2], a competition where participants are given sample codes of malware and compete by bypassing antivirus products through modifying the sample codes but without changing their function.

In this threat environment, every security vendor must constantly update their products to keep up. Given this is the case, how can a regular user *know* that their vendor is providing adequate protection against the latest threats? The obvious answer is that users should check published benchmarks; unfortunately, according to those tests, virtually every major product appears to be equivalent—they all “pass” or catch virtually all tested threats.

Researchers and industry members have recognized the reliability problem of the current testing. For example, the International Antivirus Testing Workshop [6], held in 2007, aimed to identify the problems in state-of-the-art anti-virus testing. The problems identified included testing of heuristics, malware collection and maintenance, white-list databases for testing false positives, the importance of testing in realistic environments, and questioning whether update-response-time is a good metric.

Despite such recognition, there has been considerable resistance to efforts to change the testing of commercial security products. For example, malware collection is an important part of testing the detection capabilities of anti-malware products. There is no consensus on how such collections should be curated, even though the composition of such a collection can have dramatic effects on observed detection performance. A lack of consensus, however, does not mean that there are not unspoken rules regarding such collections. In 2006, rather than collect viruses that exist “in the wild,” Consumer Reports created 5000 new viruses for testing. More than 100 security people and

executives from companies like Microsoft argued that with so many viruses already in circulation, it is not necessary to create new ones, and it brings the risk that the new viruses could be leaked into the wild [9].

While there are certainly ethical issues involved with creating new computer viruses, we believe there is a more fundamental issue: if you create malware from scratch for testing purposes, how do you know you’ve created the right kinds? In other words, how will you determine whether detection performance on synthetic test cases will correlate with performance on malware observed in practice? This issue is just one part of a much larger issue: how can you take into account all of the factors—detection mechanisms, relative frequencies of different kinds of malware, user behavior, host and network environment, changing attacker strategies and goals—that affect a product’s real world performance in a set of standardized lab tests?

We believe the simple answer is that you can’t—the task is impossible. There are simply too many variables. Researchers and companies will continue to argue about proper lab testing procedures because there is no single right answer: every test incorporates assumptions about the real world, and these assumptions cannot be evaluated in a laboratory setting.

Is there a way beyond this impasse? Perhaps, but only if we can test security technologies “in the field”—in the contexts in which they are used. Of course, such testing would involve attempting to protect real users from real threats while measuring relative performance. This approach is technically difficult, expensive, ethically challenging, and potentially very risky. We believe, however, that such testing is feasible based on experiences from the field of medicine, in the form of clinical trials.

### 3 Medical Clinical Trials

While computers and humans are very different systems, the medical field has long faced evaluation problems analogous to that of computer security. Specifically, before the 20th century there existed many potential “defenses”—treatments that promised to ensure or repair health—but people continued to be attacked and compromised (suffer and die prematurely from disease). While modern medicine has a variety of limitations, current medical practice has treatments that can reliably prevent or cure many conditions that before were debilitating or even fatal. What is remarkable about these treatments is that, in general, we don’t understand how they work: our understanding of living systems is still primitive in many ways. Despite this lack of knowledge, however, we are now able to differentiate treatments that work from those that do not. The key methodology for drawing such conclusions is the clinical trial [3].

The key insight behind clinical trials is that when studying systems (such as the human body) that are complicated, diverse, and tightly coupled with a dynamic environment, individual variables cannot be isolated and so cause and effect relationships cannot be inferred from individual observations: correlations can occur without causation, and observed effects can originate from unidentified causes. Clinical trials are an experimental methodology designed to identify causal relationships in the face of such complexity.

In medicine, clinical trials, or randomized control trials (RCTs), are planned experiments that are designed to compare treatments for a given medical condition. They use results based on a limited sample of patients to make inferences about how treatments should be conducted in the general population of patients. While the majority of clinical trials are concerned with evaluating drugs,

they can also be used to evaluate other interventions such as surgical procedures, radiotherapy, physical therapy, and diets.

To account for variations in genetic makeup, lifestyle, life history, and environment, clinical trials are designed with several key features:

**Selected populations** At risk or afflicted individuals are studied, rather than the general population.

**Extended duration** Experiments are performed for months or, ideally, years in order to evaluate longer term effects.

**Random samples** Subjects are randomly recruited from the selected population.

**Comparable Treatments** Subjects are given one of a small selection of treatments, each of which is intended to have similar effects (i.e., they treat the same condition).

**Randomly Chosen Treatments** Subjects or doctors do not choose their treatment; instead, the treatment is randomly assigned.

**Control Groups** Some subjects do not receive any treatment or are given a placebo (e.g., a sugar pill).

**Blinding** In a single blind study, subjects do not know which treatment they are receiving. In a double-blind study, the treating doctors do not know either.

**Indicators** Often the condition studied evolves over a long period of time. Rather than wait until the end (e.g., wait until the subject is cured or dead), progress is measured by observing indicators that are known to correlate with the final outcome. For example, insulin and blood sugar levels of diabetes patients are monitored in diabetes-related trials.

Due to the constraints of particular experiments, not all clinical trials will include all of these features; the more that are used, however, the greater the statistical power of the results. In other words, each of these mechanisms help with determining causal relationships. The fewer that are used, the more likely the study will only show correlation, not causation.

While clinical trials are very powerful tools for determining cause-effect relationships, they are not able to tell *why* those relationships exist. Clinical trials do not themselves provide explanations or models; what they can do, however, is test the validity and completeness of models. For example, in medicine drugs that work well in lab experiments routinely fail to work in clinical trials on people. This failure happens even when the precise molecular mechanism of the drug is known. Quite simply, we cannot capture the full complexity of the human body in any current model or lab. With clinical trials, however, we can at least make sure that only proven treatments are given to regular patients.

## 4 Computer Security Clinical Trials

Because computers are engineered systems, we are much better able to determine cause and effect in computer security than in medicine. However, while it is relatively straightforward to understand

a given vulnerability and devise a patch that fixes it, as we explained in Section 2, it is not nearly so easy to determine what produce the ultimate result of more secure systems. So, here we ask, is it potentially feasible to adapt the clinical trial methodology to computer security?

The key constraint to the feasibility question is to realize that clinical trials cannot be used to address the same questions as standard security evaluation techniques. We cannot use a clinical trial to analyze malware, expose a new software vulnerability, or test a new cryptographic protocol. However, we can use clinical trials to address questions such as the following:

- What is the security benefit of running an antivirus program on a personal computer in a typical home?
- Do personal firewalls provide additional protection for technically advanced users on their home machines?
- Does user training protect organizations from social engineering attacks?

Note the key feature of these questions is that they address interactions between computers and people. Additionally, they identify specific demographics and operational contexts.

Now, in medicine it takes a team of people to develop a clinical trial design: experts in the specific treatment must work with general clinicians, statisticians, experts in patient recruitment, ethicists, and others. Given that computer security clinical trials will also deal with human populations (along with computer populations), many of the same technical, legal, ethical, and logistical issues will need to be addressed. For these reasons, we cannot hope to present a complete trial design here; however, we can give an outline for a plausible computer security clinical trial. Here we present a sketch of a trial addressing the first question: the benefit of antivirus programs.

It is generally recommended that all personal computer users (at least, those running a version of Microsoft Windows) run an up-to-date antivirus scanner. A clinical trial designed to test their relative benefit could have the following characteristics:

**Population** Users running Microsoft Windows Vista SP2 on a home machine connected to the Internet via a large home internet service provider (ISP).

**Duration** Three years, with preliminary results reported after each year.

**Sample** 1000 ISP subscribers would be randomly recruited to participate in the trial. Each subscriber would be given the following incentives to participate: free technical support and automatic offsite backups for all machines enrolled in the trial and their users. In return, they would have to agree to researchers monitoring their computer usage (subject to appropriate privacy and other controls). Users would be allowed to drop out of the trial at any time.

**Treatments** Three major antivirus programs would be selected for the trial and randomly assigned to different households. Note that only antivirus programs would be allowed to be installed; otherwise, only the standard security software that comes with Windows Vista would be allowed to be used.

**Control** A control group would receive no antivirus program and would be prohibited from running any host-based antivirus program. However, their network communications, disks, and RAM would be scanned for malware using special hardware and/or virtual machine technologies.

**Blinding** The antivirus programs would be modified to remove any obvious corporate insignia or other advertising. Color schemes would also be modified to make them as similar as possible. Otherwise, however, their interfaces would remain the same.

In addition, the control group computers would run a program that mimicked the appearance and behavior of an antivirus program. It would provide a Windows tray icon and it periodically would report that its signatures were updated. In addition, it would check and report a variety of relatively innocuous, common problems such as tracking cookies. This program would do no proper scanning and it would provide no protection from malware.

**Indicators** A variety of measures would be required to monitor the characteristics of the monitored computers. Offsite backups would be regularly scanned for examples of known malware using a large number of commercial scanners (including ones not part of the test). A custom kernel driver would regularly report system CPU, disk, and network usage. Individuals would self-report problems to the supplied technical support service. And, a small subset of machines, say 100, would be periodically inspected by technicians to evaluate computer health and other characteristics.

While there are a variety of logistical, technological, and financial challenges implicit in the above description, it should be clear that it would be possible to run this trial given the right resources. While we could speculate on what results we might find from such a study, the fact is that we don't know what would be found. Indeed, that is the *key point* of clinical trials: they can reveal interactions and behaviors that are not observed in laboratories nor predicted by theoretical models. The most surprising result, we think, would be that we would learn nothing of consequence from such an exercise.

## 5 Objections

There are many potential objections to the use of the clinical trial methodology in a computer security context. Here we address some of the ones that have arisen in our discussions.

### 5.1 Biology vs. Computers

One significant objection is that computer security is fundamentally different from medicine because the adversaries we face are not microorganisms, but people—intelligent, motivated people. While others have long debated the merits of the biological metaphor for computer security [7], we believe that debate is not relevant to the question of computer security clinical trials. It is true that, like any experimental methodology, clinical trials are backwards looking; thus, it is always possible that their results could be invalidated via attacker innovation. Indeed, this possibility is taken into account in the design of modern security software through the use of automated update mechanisms; thus, clinical trials of security software will, implicitly, be testing the software *and*

*the organization behind it.* In practice, then, we would really be comparing humans (attackers) versus humans (defenders), as mediated by a computational battle field.

But even if we are talking about human institutions, as with many financial products, past performance is not indicative of future results. Given that we cannot predict the future of security technologies using any current technique (including formal models), however, past performance is all we have to go on when choosing security solutions. Clinical trials are merely a formal methodology for rigorously assessing that past performance.

## 5.2 Utility

Even if adopted, a clinical trial methodology will not be a panacea with respect to security. While the approach should demonstrate the real world effectiveness of products, it will not explain *why* these differences exist. For example, consider two virus scanners. Our trial would perhaps show that one product provides statistically better protection than the other—but it would not (directly) provide any explanation for their differential performance. Is it the actual rate of virus detection? The speed or ease of update? While individual users may be able to say what they liked about the product they were given, such opinions only provide clues as to the cause. As such, the results produced by the trial may be both unexpected and, *prima facie*, inexplicable. As such, a trial would show what works and what does not, but does not in any way suggest new solutions or guide new developments.

Because of these limitations, clinical trials should be seen as a complement to, not a replacement of, lab testing of security technologies. Indeed, we believe better methodologies are needed for lab evaluations as well. Our purpose here, though, is to point out that lab testing cannot be expected to address all of the issues that arise in deploying security solutions. In fact, clinical trials provide a rigorous way to determine to what extent solutions developed in the lab are applicable to practice.

## 5.3 Expense

To be sure, clinical trials are an expensive way to evaluate systems. The fundamental cost is one of labor: the subjects in the trial must be recruited, monitored, and supported by people. While some of this cost can be mitigated in the computer security context through appropriate automation, a clinical trial will always be at least an order of magnitude more expensive than a simple lab comparison.

While such expense is significant, it is also true that computer security is a multi-billion dollar market; a very small percentage of this would be enough to support many clinical trials. Further, the cost is justified by the importance of the industry. Organizations are now being required by regulation to implement security solutions. Such implementations can be very expensive. To date, we have no way of determining whether those solutions provide concrete benefits in practice.

If clinical trials are shown to work for computer security, it is likely they will become mandated by regulation, much as they have been for medicine. We think such a change would actually be to the benefit of the computer security industry. Before medical practice was regulated, there was a vigorous but relatively small trade in patent medicines—unregulated preparations that claimed to cure people's ills. Despite being pioneers in marketing and advertising, patent medicines were widely maligned and mistrusted, largely because in general they didn't actually work [8]. In



contrast, modern medicine is an extremely large, lucrative, and well-respected enterprise. If our community can, as a group, recommend solutions for which we have scientific evidence of their efficacy, perhaps computer security will also see a transformation in terms of its scope and prestige.

## 6 Conclusion

In order for the field of computer security to progress, the issue of security metrics is one which must be addressed. To this end, we have proposed applying the proven techniques used in medical clinical trials to security. This approach is attractive as it measures the performance of solutions in practice. Given the importance of information assurance in the modern world and the increasing regulatory requirements for operational security, we believe the cost and complexity of clinical trials are justified. While the ultimate value of security clinical trials will only be known in retrospect, we are optimistic that clinical trials will help the development and deployment of effective security technologies.

## 7 Acknowledgments

The authors wish to thank Tim Furlong for first thinking of the computer security clinical trial in a lab brainstorming meeting in the summer of 2006. AS, YL, and HI acknowledge support from Canada's NSERC, though the Discover Grants program and the ISSNet Strategic Network, and MITACS.

## References

- [1] Symantec Corporation. Symantec global internet security threat report, volume xiv. <http://www.symantec.com/business/theme.jsp?themeid=threatreport>, 2009.
- [2] Defcon. Race-to-zero. <http://www.racetozero.net/>.
- [3] Lawrence M. Friedman, Curt D. Furberg, and David L. DeMets. *Fundamentals of clinical trials (Third Edition)*. Springer-Verlag New York Inc., 1998.
- [4] NetworkWorld. Storm worm. <http://www.networkworld.com/news/2007/080207-black-hat-storm-worms-virulence.html>.
- [5] B. Schneier. Bruce schneier's hints for identifying crypto snake oil. <http://www.schneier.com/crypto-gram-9902.html#snakeoil>.
- [6] Web Site. International antivirus testing workshop. <http://www.f-prot.com/workshop2007/>.
- [7] Anil Somayaji, Michael Locasto, and Jan Feyereisl. Panel: The future of biologically-inspired security: Is there anything left to learn? In *Proceedings of the 2007 Workshop on New Security*, New York, NY, 2008. Association for Computing Machinery.

- [8] John Styles. Product innovation in early modern london. *Past & Present*, 168(1):124–169, 2000. Discusses how proprietary medicines were innovators in marketing but mentions that they weren’t trusted, were regarded as quack medicines, and led to general mistrust of marketing and advertisement.
- [9] WashingtonPost.com. Anti-virus testing and consumer reports. [http://voices.washingtonpost.com/securityfix/2006/08/antivirus\\_testing\\_and\\_consumer\\_1.html](http://voices.washingtonpost.com/securityfix/2006/08/antivirus_testing_and_consumer_1.html).