

Modeling the Spread of Active Worms

Zesheng Chen

Dept. of Electrical &
Computer Engineering
Georgia Institute of Technology
Atlanta, GA 30332
Email: zchen@ece.gatech.edu

Lixin Gao

Dept. of Electrical &
Computer Engineering
Univ. of Massachusetts
Amherst, MA 01002
Email: lgao@ecs.umass.edu

Kevin Kwiat

Air Force Research Lab
Information Directorate
525 Brooks Road
Rome, NY 13441-4505
Email: kwiatk@rl.af.mil

Abstract—Active worms spread in an automated fashion and can flood the Internet in a very short time. Modeling the spread of active worms can help us understand how active worms spread, and how we can monitor and defend against the propagation of worms effectively. In this paper, we present a mathematical model, referred to as the Analytical Active Worm Propagation (AAWP) model, which characterizes the propagation of worms that employ random scanning. We compare our model with the Epidemiological model and Weaver’s simulator. Our results show that our model can characterize the spread of worms effectively. Taking the Code Red v2 worm as an example, we give a quantitative analysis for monitoring, detecting and defending against worms. Furthermore, we extend our AAWP model to understand the spread of worms that employ local subnet scanning. To the best of our knowledge, there is no model for the spread of a worm that employs the localized scanning strategy and we believe that this is the first attempt on understanding local subnet scanning quantitatively.

Index Terms—security, worm, modeling

I. INTRODUCTION

Active worms have been a persistent security threat on the Internet since the Morris worm arose in 1988. The Code Red and Nimda worms infected hundreds of thousands of systems, and cost both the public and private sectors millions of dollars [1], [2], [3], [4]. Active worms propagate by infecting computer systems and by using infected computers to spread the worms in an automated fashion. Staniford et al. show that active worms can potentially spread across the Internet within seconds [5]. It is therefore of great importance to characterize and monitor the spread of active worms, and be able to derive methods to effectively defend our systems against them.

About ten years ago, Kephart and White presented the Epidemiological model to understand and control the prevalence of viruses [6], [7], [8]. This model is based on biological epidemiology and uses nonlinear differential equations to provide a qualitative understanding of virus spreading. White pointed out, however, that the “mystery” of the Epidemiological model

Zesheng Chen was with the Department of Electrical & Computer Engineering, University of Massachusetts, Amherst, MA 01002 when this work was performed.

This work is supported in part by NSF grant ANI-9977555, ANI-0085848, NSF CAREER Award grant ANI-9875513, and Air Force Research Lab. Any opinions, findings, and conclusions or recommendations expressed in this material are those of the authors and do not necessarily reflect the views of the National Science Foundation, and Air Force Research Lab.

is that it fails to predict that virtually most viruses will be slow in global prevalence [9].

In this paper, we present a model, referred to as the Analytical Active Worm Propagation (AAWP) model, which characterizes the propagation of worms that employ random scanning. We take advantage of a discrete time model and deterministic approximation to describe the spread of active worms. Our model captures the characteristics of the spread of active worms and explains the aforementioned “mystery” to some extent. In order to evaluate our model, we compare it to the simulator in [10]. Experimental results show that our model can effectively characterize the propagation of worms.

In addition to modeling the spread of worms, we attempt to answer the following questions:

- How can we monitor the spread of active worms accurately? When the Code Red v2 worm broke out on July 19th, 2001, CAIDA used one /8 network and two /16 networks to monitor the spread [11]. It is not clear, however, whether the data collected from these networks can reflect the actual spread of the worm. If the data does not reflect the actual spread of the worm, what is the size of the network that should be used to monitor the infected machines? Our results show that monitoring a /8 network is sufficient for characterizing the spread of active worms accurately.
- How can we detect the spread of active worms in a timely fashion? To the best of our knowledge, no effective worm detection mechanism is available. One simple detection system uses unused IP addresses to detect the scans from active worms. With the help of the AAWP model, we derive the number of IP addresses needed for detecting the spread of active worms effectively. Although this simple detection system might generate false alarms, we believe that it is the first step in understanding the effectiveness of the detection system quantitatively.
- How can we defend against the spread of active worms effectively? We perform a study on how well a worm defending tool can slow down the propagation of worms based on our model. Our study quantitatively illustrates the size of address space needed to stop or slow down the Code Red v2 like worms effectively.

Furthermore, developing an analytical model for the spread

of a worm employing a localized scanning strategy is significantly more difficult than that for random scanning [5]. We extend the AAWP model to characterize the spread of a worm that employs the localized scanning strategy, which is used by the Code Red II and Nimda worms. Our model shows that worms that employ localized scanning spread at a slower rate than those employing random scanning despite the fact that localized scanning can potentially penetrate beyond firewalls. To the best of our knowledge, this is the first attempt in understanding the local subnet scanning policy quantitatively.

The remainder of this paper is structured as follows. Section II describes how active worms spread, and introduces the parameters for characterizing their propagation. In Section III, we present the AAWP model, and compare it to the Epidemiological model and Weaver's simulator. In addition, we use the AAWP model to simulate the spread of the Code Red v2 worm. Section IV outlines the applications of the AAWP model, which include verifying the correctness of the worm's monitoring data, developing a detection system and evaluating the LaBrea defense system. In Section V, we extend the AAWP model to worms that employ local subnet scanning. We conclude this paper in Section VI with a brief summary and an outline of our future work.

II. SPREAD OF ACTIVE WORMS

In this section, we first describe how active worms spread, then introduce the parameters used in the spread of active worms. Finally, we present two worm scanning models: random scanning and local subnet scanning.

When an active worm is fired into the Internet, it simultaneously scans many machines in an attempt to find a vulnerable machine to infect. When it finally finds its prey, it sends out a probe to infect the target. If successful, a copy of this worm is transferred to this new host. This new host then begins running the worm and tries to infect other machines. When an invulnerable machine or an unused IP address is reached, the worm poses no threat. During the worm's spreading process, some machines might stop functioning properly, forcing the users to reboot these computers or at least kill some of the processes that may have been exploited by the worm. Then these infected machines become vulnerable machines again, and are still inclined to further infection. When the worm is detected, people will try to slow it down or stop it. A patch, which repairs the security hole of the machines, is used to defend against worms. When an infected or vulnerable machine is patched, it becomes an invulnerable machine.

To speed up the spread of active worms, Weaver presented the "hitlist" idea [10]. Long before an attacker releases the worm, he/she gathers a list of potentially vulnerable machines with good network connections. After the worm has been fired onto an initial machine on this list, it begins scanning down the list. Hence, the worm will first start infecting the machines on this list. Once this list has been exhausted, the worm will then start infecting other vulnerable machines. The machines on this list are referred to as the "hitlist". After the worm infects

the hitlist rapidly, it uses these infected machines as "stepping stones" to search for other vulnerable machines. In this paper we do not consider the amount of time it takes a worm to infect the hitlist since the hitlist can be acquired well before a worm is released and be infected in a very short period of time. Table I shows the parameters involved in the spread of active worms.

There are several different scanning mechanisms that active worms employ, such as random, local subnet, permutation and topological scanning [5]. In this paper we focus on two mechanisms, random scanning and local subnet scanning. In random scanning, it is assumed that every computer in the Internet is just as likely to infect or be infected by other computers. Such a network can be pictured as a fully-connected graph in which the nodes represent computers and the arcs represent connections (neighboring-relationships) between pairs of nodes. This topology is called "homogeneous mixing" in the theoretical epidemiology [7]. Our AAWP model is used to model random scans. In local subnet scanning, computers also connect to each other directly, forming "homogeneous mixing". However, instead of selecting targets randomly, the worms preferentially scan for hosts on the "local" address space. For example, the Nimda worm selects target IP addresses as follows [3]:

- 50% of the time, an address with the same first two octets will be chosen.
- 25% of the time, an address with the same first octet will be chosen.
- 25% of the time, a random address will be chosen.

We will extend the AAWP model to the Local AAWP (LAAWP) model in Section V to understand the function of the propagation parameters and analyze the spread of active worms that employ local subnet scanning.

III. MODELING THE SPREAD OF ACTIVE WORMS THAT EMPLOY RANDOM SCANNING

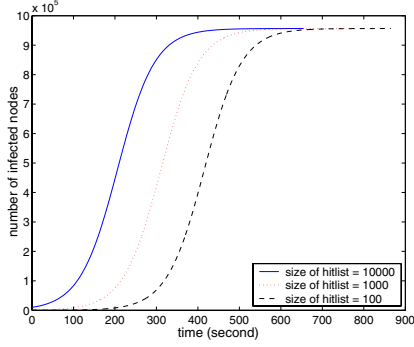
To understand the characteristics of the spread of active worms that employ random scanning, we develop the AAWP model, which uses the discrete time and continuous state deterministic approximation model. In this section, we first describe the AAWP model in detail, then compare it to the Epidemiological model and Weaver's simulator, finally use it to simulate the Code Red v2 worm.

A. Deterministic Approximation Modeling

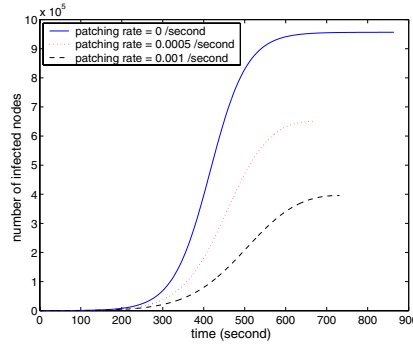
We assume that worms can simultaneously scan many machines and will not re-infect a machine that is already infected. We also assume that the machines on the hitlist are already infected at the start of the worm's propagation. Suppose that an active worm takes one time tick to complete infection. That is, when one scan hits a machine, regardless of whether this machine is vulnerable, invulnerable, infected or with an unused IP address, the time it takes for the worm to finish communicating with this machine is one time tick. This assumption might not be realistic, but it can simplify the model without significantly affecting the results.

TABLE I
THE PARAMETERS FOR THE SPREAD OF ACTIVE WORMS

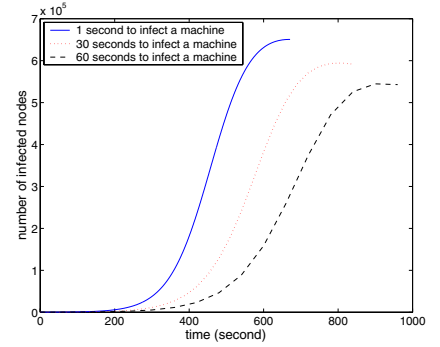
Parameters	Notation	Explanation
# of vulnerable machines	N	the number of vulnerable machines
Size of hitlist	h	the number of infected machines at the beginning of the spread of active worms
Scanning rate	s	the average number of machines scanned by an infected machine per unit time
Death rate	d	the rate at which an infection is detected on a machine and eliminated without patching
Patching rate	p	the rate at which an infected or vulnerable machine becomes invulnerable



(a) Effect of Hitlist Size (All cases are without patching and take a period of one second to complete infection.)



(b) Effect of Patching Rate (All cases have a hitlist of 100 entries and take a period of one second to complete infection.)



(c) Effect of Time to Complete Infection (All cases have a hitlist of 100 entries and a patching rate of 0.0005 /second.)

Fig. 1. Modeling the Spread of Active Worms that Employ Random Scanning (All cases are for 1,000,000 vulnerable machines, a scanning rate of 100 scans/second, and a death rate of 0.001 /second.)

Although the Internet's address space isn't completely connected, active worms always scan 2^{32} entry addresses. Therefore, for random scanning, the probability that any computer is hit by one scan is $\frac{1}{2^{32}}$. Let m_i and n_i denote the total number of vulnerable machines (including the infected ones) and the number of infected machines at time tick i ($i \geq 0$) respectively. Before the active worms spread ($i = 0$), $m_0 = N$ and $n_0 = h$.

Theorem 1: If there are m_i vulnerable machines (including the infected ones), and n_i infected computers, then on average, the next time tick will have $(m_i - n_i)[1 - (1 - \frac{1}{2^{32}})^{sn_i}]$ newly infected machines, where s is the scanning rate.

PROOF: Let e_i denote the number of newly infected machines at time tick i ($i \geq 0$). n_i infected machines can generate sn_i scans in an attempt to infect other machines. So if we can prove that $E\{e_{i+1}/k\} = (m_i - n_i)[1 - (1 - \frac{1}{2^{32}})^k]$ for any k ($k > 0$) scans, then the equation also holds when $k = sn_i$.

We prove the above equation by induction on k . When $k = 1$, since there are $(m_i - n_i)$ vulnerable machines that have not yet been infected, the probability that one scan can add a newly infected machine is $\frac{m_i - n_i}{2^{32}}$, which is equivalent to $(m_i - n_i)[1 - (1 - \frac{1}{2^{32}})^1]$. Suppose that the theorem is true for $k = j$, i.e., $E\{e_{i+1}/k = j\} = (m_i - n_i)[1 - (1 - \frac{1}{2^{32}})^j]$. Then, when $k = j + 1$, we divide $j + 1$ scans into two parts: the first j scans and the last scan. There are two possibilities for the last scan: adding a newly infected machine or not. Let

the variable $Y = 1$ if the last scan hits a vulnerable machine that has not yet been infected and let $Y = 0$ otherwise. Then,

$$\begin{aligned}
 & E\{e_{i+1}/k = j + 1\} \\
 &= (E\{e_{i+1}/k = j\} + 1)P(Y = 1) + E\{e_{i+1}/k = j\} \cdot P(Y = 0) \\
 &= (E\{e_{i+1}/k = j\} + 1) \frac{m_i - n_i - E\{e_{i+1}/k = j\}}{2^{32}} + \\
 & E\{e_{i+1}/k = j\} (1 - \frac{m_i - n_i - E\{e_{i+1}/k = j\}}{2^{32}}) \\
 &= \frac{m_i - n_i}{2^{32}} + (1 - \frac{1}{2^{32}}) E\{e_{i+1}/k = j\} \\
 &= (m_i - n_i) [1 - (1 - \frac{1}{2^{32}})^{j+1}]
 \end{aligned}$$

which means that when $k = j + 1$, it is also true. Therefore, when $k = sn_i$, $E\{e_{i+1}/k = sn_i\} = (m_i - n_i)[1 - (1 - \frac{1}{2^{32}})^{sn_i}]$. That is, on the next time tick there will be $(m_i - n_i)[1 - (1 - \frac{1}{2^{32}})^{sn_i}]$ expected newly infected machines. ■

Given death rate d and patching rate p , on the next time tick there will be $dn_i + pn_i$ infected machines that will change to either vulnerable machines without being infected or invulnerable machines, and the total number of vulnerable machines (including the infected ones) will be reduced to $(1 - p)m_i$. Therefore, on the next time tick the number of total infected

machines will be $n_{i+1} = n_i + (m_i - n_i)[1 - (1 - \frac{1}{2^{32}})^{sn_i}] - (d+p)n_i$. At the same time, $m_{i+1} = (1-p)m_i$, which gives $m_i = (1-p)^i m_0 = (1-p)^i N$. That is,

$$n_{i+1} = (1-d-p)n_i + [(1-p)^i N - n_i][1 - (1 - \frac{1}{2^{32}})^{sn_i}] \quad (1)$$

where $i \geq 0$ and $n_0 = h$. The recursion process will stop when there are no more vulnerable machines left or when the worm cannot increase the total number of infected machines.

Using Equation (1), we can find the characteristics of the active worms' spreading. For example, Figure 1(a) shows the propagation of the active worms with different hitlist sizes. As the size of the hitlist increases, it takes the worms less time to spread. Figure 1(b) depicts another example. As the patching rate grows, the spread of active worms slows down. This complies with our intuition. It should be noted that because the patching rate $p > 0$, the two slower curves return to zero at the end. Here, we only draw the uprising part of curves and ignore the falling part.

At the beginning, we assume that it takes the worms one time tick to infect a machine. To display the effect of the amount of time it takes to infect a machine on the worm propagation, we simply change the time unit. For example, in Figure 1(c) we first draw the curve with a time interval of one second, which is the amount of time required to complete infection. If the worm needs 30 seconds to infect a machine, we set the time unit to 30 seconds, and change the corresponding s, d, p parameters for this period of time. In this case, the parameters s, d, p will become $30s, 30d, 30p$ for a time period of 30 seconds. Then, we can use the AAWP model to get the result. But, now n_i expresses the number of infected machines at $30i$ seconds ($i \geq 0$). This figure shows the effect of the time to complete infection on the worm's propagation. The worm's propagation will be slowed down as the time required to infect a machine increases.

We can change the values of the parameters N, h, s, d, p and the time to complete infection in the AAWP model to observe how active worms spread. This model can be used to quantitatively explain how we can monitor the spread of active worms, develop a sensor detection system, and evaluate the LaBrea tool defense system, which we will cover later.

B. Comparing Our AAWP Model to the Epidemiological Model and Weaver's Simulator

In the Epidemiological model, a nonlinear differential equation is used to measure the virus population dynamics [7]:

$$\frac{dn}{dt} = \beta n(1-n) - dn$$

where $n(t)$ is the fraction of infected nodes, β is the birth rate (the rate at which an infected machine infects other vulnerable machines), and d is the death rate. The solution to the above equation is

$$n(t) = \frac{n_0(1-\rho)}{n_0 + (1-\rho-n_0)e^{-(\beta-d)t}} \quad (2)$$

where $\rho = \frac{d}{\beta}$ and $n_0 \equiv n(t=0) = \frac{\text{size of hitlist}}{N} = \frac{h}{N}$.

In fact, we can easily deduce the relationship between the birth rate and the scanning rate: $\beta = \frac{Ns}{2^{32}}$.

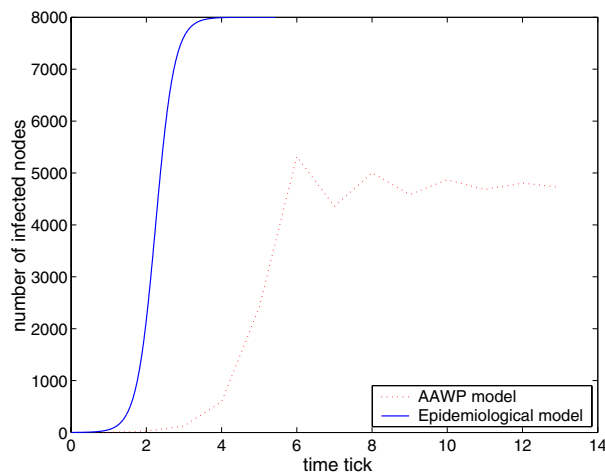
It is interesting that when the Code Red v2 worm surged in July of 2001, Staniford also independently presented the same model to explain the Random Constant Spread (or RCS) theory of the Code Red v2 worm [5]. Zou extended the Epidemiological model to the two-factor worm model, which takes consideration of the human countermeasure and the worm's impact on Internet traffic and infrastructure [12].

The differences between the AAWP model and the Epidemiological model are:

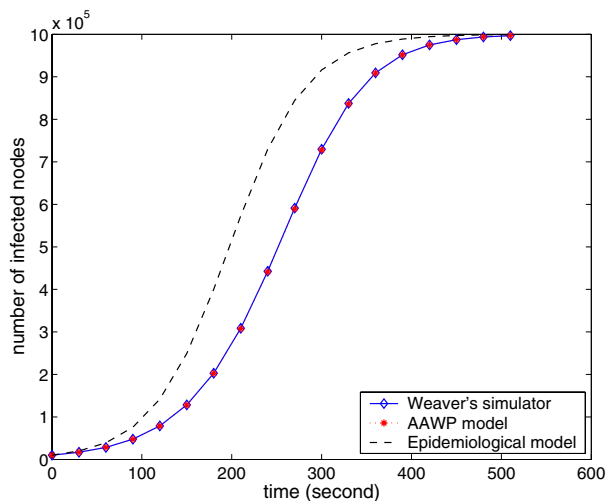
- 1) The Epidemiological model uses a continuous time differential equation, while the AAWP model is based on a discrete time model. We believe that the AAWP model is more accurate. Because in the AAWP model, a computer cannot infect other machines before it is infected completely. But in the Epidemiological model, a computer begins devoting itself to infecting other machines even though only a "small part" of it is infected. Therefore, the speed that the worm can achieve and the number of machines that can be infected are totally different.
- 2) The Epidemiological model neither considers the patching rate nor the time that it takes the worm to infect a machine, while the AAWP model does. During the propagation of the worm, it is possible nowadays to promptly patch the vulnerability on computers, assuming a reasonable patching rate. And different worms have different infection abilities which are reflected by the scanning rate (or the birth rate) and the time taken to infect a machine. The time required to infect a machine always depends on the size of the worm's copy, the degree of network congestion, the distance between source and destination, and the vulnerability that the worm exploits. From Figure 1(c), it can be seen that the time to infect a machine is an important factor for the spread of active worms.
- 3) In the AAWP model, we consider the case that the worm can infect the same destination at the same time, while the Epidemiological model ignores the case. In fact, it is not uncommon for a vulnerable machine to be hit by two (or more) scans at the same time.

Both models, however, try to get the expected number of infected machines, given the size of the hitlist, total number of vulnerable machines, scanning rate/birth rate and death rate. The Epidemiological model can easily deduce the closed form and can be used in topology orientation, such as E-mail worms or peer-to-peer worms. In this paper, we focus on active worms that select destinations randomly or employ local subnet scanning, such as the Code Red and Nimda worms. Hence, the AAWP model, which is built on the "homogeneous mixing" topology, is sufficient for our work.

Figure 2(a) shows the comparison between these two models with 10,000 vulnerable machines, a hitlist with 1 entry, a birth rate of 5 /time tick and a death rate of 1 /time tick (the parameters are from [7]). It takes the Epidemiological

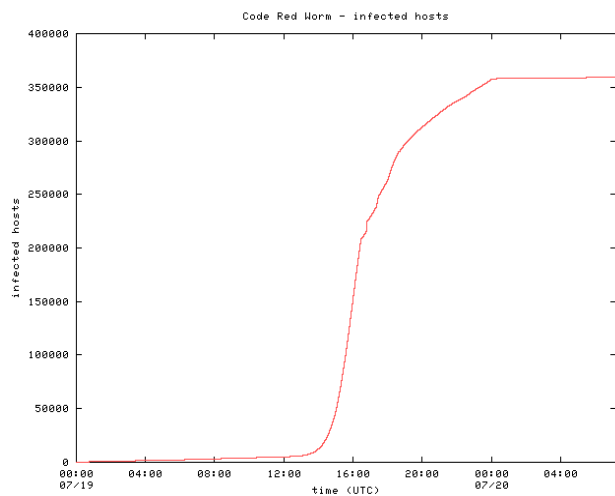


(a) All cases are for 10,000 vulnerable machines, a hitlist with 1 entry, a scanning rate of 2147500 scans/time tick or a birth rate of 5 /time tick and a death rate of 1 /time tick. No patching and a time period of 1 time tick to complete infection for the AAWP model.

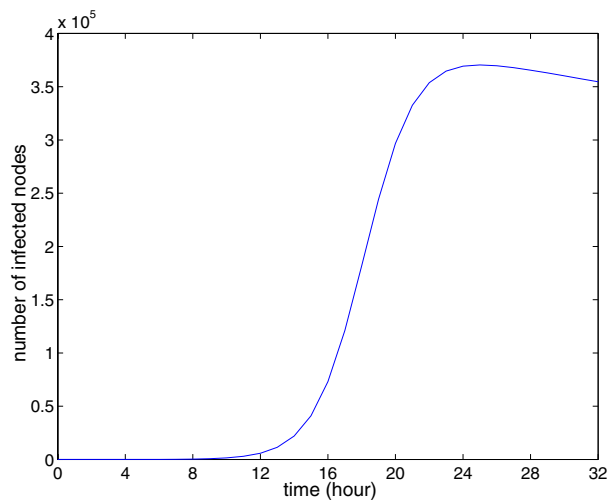


(b) All cases are for 1,000,000 vulnerable machines, a hitlist with 10,000 entries and a scanning rate of 100 scans/second. A time period of 30 seconds to complete infection for Weaver's simulator and the AAWP model. A death rate of zero for both the AAWP model and the Epidemiological model. No patching for the AAWP model.

Fig. 2. Comparing the AAWP Model to the Epidemiological model



(a) Measurement of the Code Red v2 worm spread using real data from CAIDA.



(b) A simulation of the spread of the Code Red v2 worm (500,000 vulnerable machines, starting on a single machine, a scanning rate of 2 scans/second, a death rate of 0.00002 /second, a patching rate of 0.000002 /second, and a time period of 1 second to complete infection).

Fig. 3. Real Data from CAIDA [11] and Simulated Code Red v2 Like Worm from the AAWP Model

model about 4 time ticks to enter an equilibrium stage, while the AAWP model needs about 10 time ticks. Moreover, after entering the equilibrium stage, the Epidemiological model totally infects 8,000 vulnerable machines (occupying 80% of all vulnerable machines), while the AAWP model infects about

4,750 vulnerable machines (occupying 47.5% of all vulnerable machines). This shows that our model can better explain the low level of worm prevalence in [9].

Weaver wrote a small, abstract simulator of a Warhol worm's spread [10]. This simulator uses a 32-bit, 6-round variant

of RC5 to generate all permutations and random numbers. We only modified one condition of this simulator to fit the assumption which we presented above. That is, all “newly” infected machines on a previous time tick will be activated at the same time on the current time tick, other than based on different clocks. Figure 2(b) shows the growing of infected nodes with time for the two models and Weaver’s simulator, which have the following parameters: a total of 1,000,000 vulnerable machines, a hitlist of size 10,000, a scanning rate of 100 scans/second, a death rate of zero, no patching, and a time period of 30 seconds to infect one machine. This figure shows that the AAWP model and Weaver’s simulator results overlap. While our model and Weaver’s simulator take about 6 minutes to infect 90% of the vulnerable machines, the Epidemiological model only takes about 5 minutes.

C. Simulating the Code Red v2 Worm

On July 19th, 2001, the Code Red v2 worm infected more than 359,000 computers in less than 14 hours [11]. This worm spreads by probing random IP addresses and infecting all the hosts that are vulnerable to the IIS exploit. CAIDA [13] collected real data to measure the spread of the Code Red v2 worm. The data were collected from two locations: one /8 network at UCSD and two /16 networks at Lawrence Berkeley Laboratory (LBL). In these data, hosts were considered to be infected if they sent TCP SYN packets on port 80 to nonexistent hosts on these networks. Figure 3(a) shows the number of infected hosts over time [11].

We suppose that there are 500,000 vulnerable machines in the Internet, the Code Red v2 worm starts on a single machine, it performs 2 scans per second and takes one second to infect a machine. Figure 3(b) shows the spread of the simulated Code Red v2 like worm using our AAWP model, with a death rate of 0.00002 /second and a patching rate of 0.000002 /second. Because of the patching rate, the curve goes down slightly after the worm spreads for one day.

IV. APPLICATIONS OF THE AAWP MODEL

A good model can reflect the spread of real worms and at the same time resolve many practical task. In this section, we apply the AAWP model to monitoring, detecting and defending against the spread of active worms.

A. Monitoring the Spread of Active Worms

How to monitor the spreading rate of active worms is an interesting task. It has come to our attention that CAIDA collected real data from one /8 network at UCSD and two /16 networks at LBL [11]. Can these collected data reflect the actual propagation of the Code Red v2 worm? Of course these data are only the lower bound of the spread of the Code Red v2 worm. But, how much do they deviate from the reality?

Suppose that we can collect data from 2^{32-l} ($0 \leq l \leq 32$) addresses to estimate the spread of active worms. Here, l network is the special case of 2^{32-l} addresses. These addresses

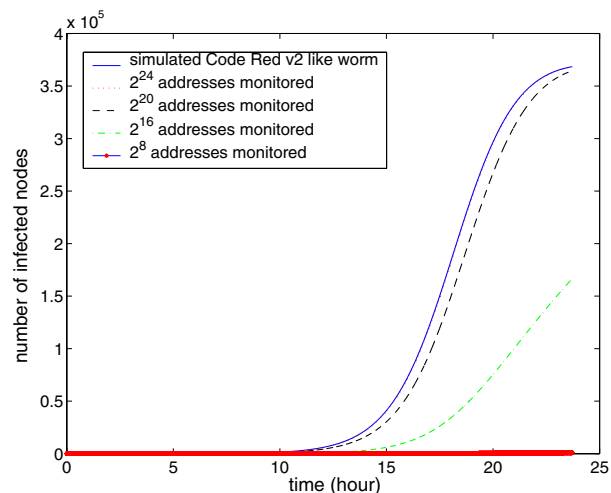


Fig. 4. Collecting data from different address spaces. All cases were for 500,000 vulnerable machines, starting on a single machine, a scanning rate of 2 scans/second, a death rate of 0.00002 /second, a patching rate of 0.000002 /second, and a time period of 1 second to complete infection.

are considered unused IP addresses. When the scans from the infected machine hit any address in this space, it is counted if and only if it has not been counted before. The probability that one scan hits this space is $\frac{2^{32-l}}{2^{32}} = \frac{1}{2^l}$. If active worms can generate s scans per time tick, then the probability that an uncounted infected machine is observed on the next time tick is $prob = 1 - (1 - \frac{2^{32-l}}{2^{32}})^s = 1 - (1 - \frac{1}{2^l})^s$. Furthermore, if $2^l \gg 1$ and $2^l \gg s$, then

$$prob = 1 - (1 - \frac{1}{2^l})^s \approx 1 - e^{-\frac{s}{2^l}} \approx \frac{s}{2^l} \quad (3)$$

Let A_i denote the number of observed infected machines at time tick i ($i \geq 0$). Before time tick $i + 1$, there are $n_i - A_i$ uncounted infected machines. The events that uncounted infected machines are observed are independent of one another. Hence, the number of “newly” observed infected machines satisfies the Binomial distribution. Then, at time tick $i + 1$ the expected number of “newly” observed infected machines is $prob \cdot (n_i - A_i)$. Therefore,

$$A_{i+1} = A_i + prob \cdot (n_i - A_i) = (1 - prob)A_i + prob \cdot n_i \quad (4)$$

where, $i \geq 0$ and $A_0 = 0$.

Based on the AAWP model, we can evaluate the effect of the different address spaces from which we collect data. Figure 4 shows one example in which we simulate the Code Red v2 worm. The curve where 2^{24} addresses are monitored is very close to the “real” worm propagation using the AAWP model. The curve where 2^{20} addresses are monitored grows at a slower rate than the curve where 2^{24} addresses are monitored, but at the same time is a much better representation than the curve where 2^{16} addresses are monitored. The curve where 2^8 addresses are monitored gives the worst results, which can be understood from Equation (3): when $l = 24$, $prob \approx 0$, then $A_{i+1} \approx A_i$, making the curve a horizontal line along the x-axis.

From the analysis above, we conclude that monitoring 2^{24} addresses gives us a better representation of the propagation of active worms. But an address space smaller than 2^{20} is not adequate to observe the actual spread of active worms.

B. Detection Speed

One of the goals of modeling the spread of active worms is to be able to detect them. Here, we present a simple and useful sensor detection system and use the AAWP model to evaluate its performance.

1) *Methodology*: It is vital to detect active worms effectively. In the near future active worms may spread across the whole Internet in a very short period of time [10], making the average detection time critical.

It is easy to figure out one simple detection system. First, put some sensors in the Internet to monitor a set of unused IP addresses. When the random scans from active worms hit these IP addresses, they are detected by the sensors. However, if the worms' designers know which unused IP addresses monitored by sensors, they could launch DoS attacks by sending many packets to the sensors, causing them to generate many false alarms. Therefore, sensors must have the intelligence to distinguish between the scans from active worms and DoS attacks, which requires a more complex sensor detection system. However, this challenge is beyond the scope of this paper.

For this simple detection system, some interesting questions need to be answered:

- How many unused IP addresses should be monitored by sensors in order to detect active worms rapidly?
- Given the number of IP addresses monitored, what is the average time required to detect worms?

2) *Performance of the Sensor Detection System*: The performance of the sensor detection system depends mainly on the detection time. An ideal detection system should be able to detect active worms at the beginning of their propagation. We use the average detection time as our performance indicator for the sensor detection system. Let T_d denote the detection time. Below, we will deduce the relationship between the average detection time and the number of unused IP addresses that are monitored.

Suppose that there are u unused IP addresses monitored by sensors. For a single scan, the probability that it is detected by sensors is $\frac{u}{2^{32}}$. Thus, for k scans, the probability that at least one scan is detected by sensors is $1 - (1 - \frac{u}{2^{32}})^k$.

Let D_i indicate the probability that a worm is detected at time tick i ($0 \leq i \leq j+1$), where $D_0 = 0$. Also note that at time tick j there are either no more vulnerable machines or the active worms cannot increase the total number of infected machines. Here, we assume that even if sensors fail to detect active worms, people will finally detect them, which means $D_{j+1} = 1$. Since n_{i-1} infected machines can generate sn_{i-1} scans,

$$D_i = 1 - (1 - \frac{u}{2^{32}})^{sn_{i-1}} \quad (5)$$

where $1 \leq i \leq j$. Then the expected value of detection time T_d is:

$$E\{T_d\} = \sum_{k=1}^{j+1} k \cdot [\prod_{l=0}^{k-1} (1 - D_l)] \cdot D_k \quad (6)$$

Based on the above formula and the AAWP model, Figure 5(a) shows the relationship between the average detection time and the number of unused IP addresses that are monitored by sensors when the active worms spread with varying hitlist sizes. From this figure, we know that in the case of a simulated Code Red v2 like worm (size of hitlist = 1), when monitoring 2^{24} addresses, the average detection time is only about two minutes; when monitoring 2^{16} addresses, the average detection time increases up to about two hours. If we want to detect this worm in one hour, more than 2^{18} unused IP addresses must be monitored by sensors. On the other hand, although the worm with a larger hitlist spreads faster, it can also be detected in a shorter period of time. For example, when monitoring 2^{16} IP addresses, we need about two hours to detect a worm starting on a single machine, while we need 36 minutes to detect a worm with a hitlist of 10 machines and only 5 minutes to detect a worm with a hitlist of 100 machines. Table II shows some sample results in Figure 5(a).

This simple detection system can be easily extended to more complex systems. For example, to reduce the number of false alarms, the sensors should receive several scans during a period of one time tick before the system can actually detect worms. Set S_n to be the least number of scans received by a sensor to generate an alarm during the period of one time tick. A larger S_n value generates less false alarms. The probability of detection, however, is also reduced. Based on the definition of S_n , the Equation (5) becomes

$$D_i = 1 - \sum_{l=0}^{S_n-1} \frac{(sn_{i-1})!}{l!(sn_{i-1}-l)!} (1 - \frac{u}{2^{32}})^{sn_{i-1}-l} (\frac{u}{2^{32}})^l$$

where $1 \leq i \leq j$. Figure 5(b) shows the effect of S_n on the sensor detection system. Even though a large S_n value can reduce the number of false alarms, it reduces the overall performance of the system. So reliability and performance are a tradeoff. Also, if a system with a large S_n value does not monitor enough addresses, it cannot completely detect the worms. For example, when $S_n = 5$ and the addresses space monitored is less than 2^{11} , the system cannot detect the worms even if the whole Internet has been infected.

C. Effectiveness of the Defense System

Another goal of modeling the spread of active worms is to defend against active worms. Here, we extend the AAWP model to analyze the LaBrea tool, which is put forward by Liston to slow down or even stop the spread of active worms [14].

1) *LaBrea*: LaBrea is a tool that takes over unused IP addresses on a network and creates "virtual machines" that answer to connection requests [15]. LaBrea replies to those connection requests in such a way that causes the machine at the other end to get "stuck". One can intentionally hold a connection open

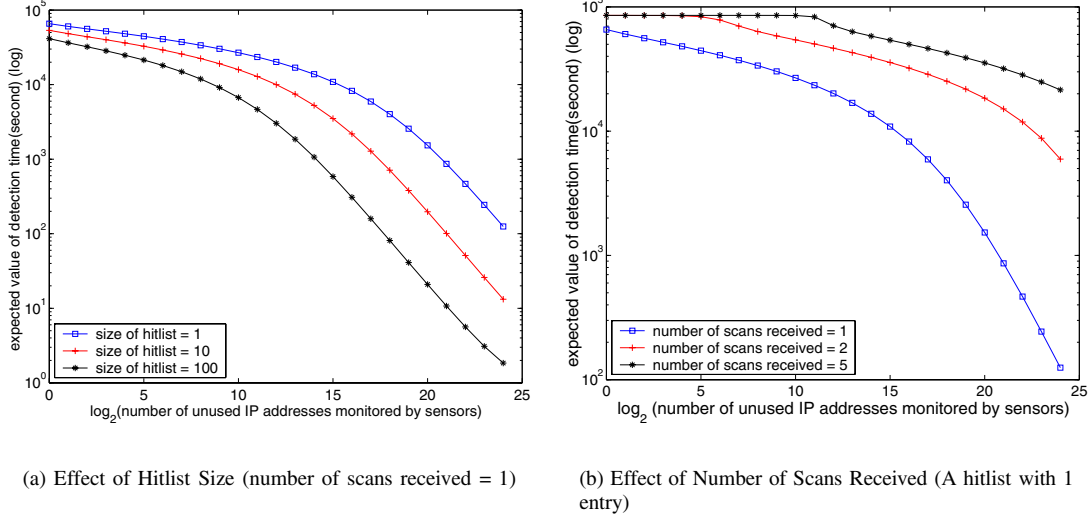


Fig. 5. Performance of sensor detection system. All cases are for 500,000 vulnerable machines, a scanning rate of 2 scans/second, a death rate of 0.00002 /second, a patching rate of 0.000002 /second, and a time period of 1 second to complete infection.

TABLE II
AVERAGE DETECTION TIME WITH HITLISTS OF DIFFERENT SIZES AND DIFFERENT NUMBERS OF UNUSED IP ADDRESSES MONITORED BY SENSORS (SECOND)

number of IP addresses monitored	2^{12}	2^{14}	2^{16}	2^{18}	2^{20}	2^{22}	2^{24}
size of hitlist = 1	20120.00	13800.00	8241.50	4021.90	1530.60	466.28	125.00
size of hitlist =10	10007.00	5267.30	2184.70	711.70	197.14	51.14	13.25
size of hitlist =100	3030.40	1065.70	308.20	81.06	20.90	5.63	1.84

for as long as he/she wishes. That is, the LaBrea tool monitors all traffic destined for some unused IP addresses. When one scan hits these IP addresses (“virtual machines”), the LaBrea tool will reply and establish a connection with the infected machine. This connection can last for a very long time.

Before we apply the LaBrea tool extensively, we should first attempt to answer one question: How many unused IP addresses should be monitored by the LaBrea tool to defend against active worms effectively?

2) *Performance of the LaBrea Tool Defense System:* Assume that the LaBrea tool is installed in the Internet and is monitoring u unused IP addresses. Suppose that now there are k scans from infected machines, beginning to search the Internet. Because the LaBrea tool can trap the scanning threads, after one time tick, there will be $\frac{u}{2^{32}}k$ scanning threads trapped. That is, there will only be $(1 - \frac{u}{2^{32}})k$ scanning threads left.

Let k_i and e_i denote the average number of scans and the number of newly infected machines at time tick i ($i \geq 0$) respectively. Taking into consideration that the LaBrea tool will affect the total number of scans, we extend Equation (1) to

$$\begin{aligned}
 m_i &= (1 - p)^i N \\
 k_{i+1} &= (1 - d - p)k_i \left(1 - \frac{u}{2^{32}}\right) + s e_i \\
 e_{i+1} &= (m_i - n_i) \left[1 - \left(1 - \frac{1}{2^{32}}\right)^{k_{i+1}}\right]
 \end{aligned}$$

$$n_{i+1} = (1 - d - p)n_i + e_{i+1}$$

where $i \geq 0$, $k_0 = 0$ and $e_0 = n_0 = h$. The recursion process will stop when there are no more vulnerable machines left or when the worm cannot increase the total number of infected machines. It should be noted that if $u = 0$, the set of formulae outlined above turn out to be the same as Equation (1).

Figure 6 shows a simulation of a Code Red v2 like worm spreading. When the LaBrea tool monitors less than 2^{16} unused IP addresses, the worm spread is only slightly affected. But when more than 2^{18} unused IP addresses are monitored, the LaBrea tool is able to effectively defend against the worm propagation. We can also see that the total number of infected machines stops increasing before all the vulnerable machines are actually infected when the LaBrea monitors more than 2^{18} unused IP addresses.

Therefore, the LaBrea tool can really slow down or stop the spread of active worms. However, at least 2^{18} unused IP addresses are needed to defend against active worms effectively. It might not be easy to persuade many network administrators to install the LaBrea. If we can get one unused class A subnet (an address space of 2^{24} addresses), which is not publicly advertised, and install the LaBrea tool to monitor the traffic into this subnet, this seems to be a good start for fighting against active worms.

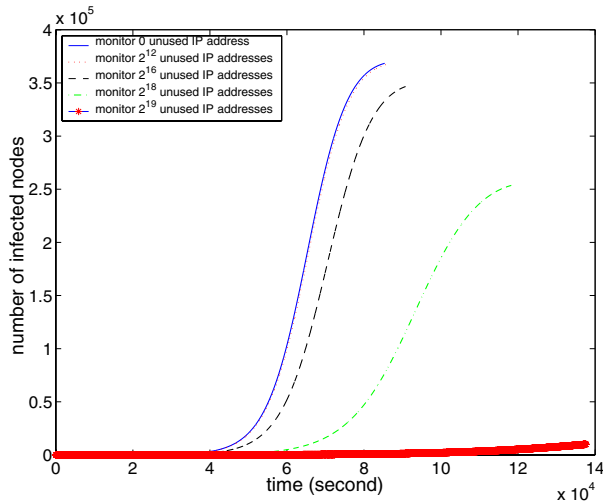


Fig. 6. Performance of the LaBrea tool detection system. All cases are for 500,000 vulnerable machines, starting on a single machine, a scanning rate of 2 scans/second, a death rate of 0.00002 /second, a patching rate of 0.000002 /second, and a time period of 1 second to complete infection.

V. MODELING THE SPREAD OF ACTIVE WORMS THAT EMPLOY LOCAL SUBNET SCANNING

Instead of simply selecting destinations at random, the Code Red II and the Nimda worms preferentially search for targets on the “local” address space [1], [3]. In this part, we extend the AAWP to the Local AAWP (LAAWP) model to understand the characteristics of the spread of active worms that employ local subnet scanning.

A. LAAWP Model

As the AAWP model, the LAAWP model uses deterministic approximation. We focus on the active worms’ scanning policy and ignore both the death rate and the patching rate to simplify the model. The function of firewalls is not considered, either.

Now suppose that a worm scans the Internet as follows:

- p_0 of the time, a random address will be chosen
- p_1 of the time, an address with the same first octet will be chosen
- p_2 of the time, an address with the same first two octets will be chosen

where, $p_0 + p_1 + p_2 = 1$. We can regard random scanning as one special case of local subnet scanning, when $p_0 = 1$, $p_1 = 0$, and $p_2 = 0$.

Assume that the vulnerable machines are evenly distributed in every subnet which is identified by the first two octets. The subnets can be classified into three different kinds of networks:

- A “special” subnet (denoted by Subnet type 1), which always has a larger hitlist size.
- $2^8 - 1$ subnets having the same first octet as the “special” subnet (denoted by Subnet type 2).
- Other $2^{16} - 2^8$ subnets (denoted by Subnet type 3).

Different kinds of networks have hitlists of different sizes. In the same type of subnet, all networks have the same hitlist size.

Let h_1 , h_2 , and h_3 denote the size of the hitlist in Subnet type 1, 2, and 3, respectively.

Let b_1 , b_2 , and b_3 denote the average number of infected machines in Subnet type 1, 2, and 3, respectively. And let k_1 , k_2 , and k_3 denote the average number of scans hitting Subnet type 1, 2, and 3, respectively. Then at some time tick, the relationship between the average number of scans hitting Subnet type i ($i = 1, 2, \text{ or } 3$) and the average number of infected machines in different Subnets is

$$\begin{aligned} k_1 &= p_2 s b_1 + p_1 s [b_1 + (2^8 - 1) \cdot b_2] / 2^8 + \\ &\quad p_0 s [b_1 + (2^8 - 1) \cdot b_2 + (2^{16} - 2^8) \cdot b_3] / 2^{16} \\ k_2 &= p_2 s b_2 + p_1 s [b_1 + (2^8 - 1) \cdot b_2] / 2^8 + \\ &\quad p_0 s [b_1 + (2^8 - 1) \cdot b_2 + (2^{16} - 2^8) \cdot b_3] / 2^{16} \\ k_3 &= p_2 s b_3 + p_1 s b_3 + \\ &\quad p_0 s [b_1 + (2^8 - 1) \cdot b_2 + (2^{16} - 2^8) \cdot b_3] / 2^{16} \end{aligned}$$

For k_i ($i = 1, 2, \text{ or } 3$), the first item is the average number of scans coming from the local subnet (with the same first two octets). The second item is the average number of scans coming from neighboring subnets (with the same first octet). And the last item is the average number of scans coming from global subnets.

In every subnet the scans will randomly hit targets, which can be modeled by the AAWP model. The total number of machines will be 2^{16} , instead of 2^{32} , and the total number of scans will be k_i . Thus, Equation (1) becomes

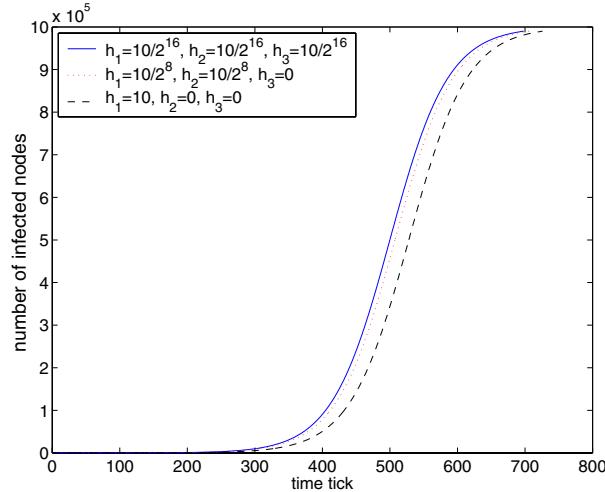
$$b'_i = b_i + \left(\frac{N}{2^{16}} - b_i \right) \left[1 - \left(1 - \frac{1}{2^{16}} \right)^{k_i} \right] \quad (7)$$

where, $i = 1, 2, \text{ or } 3$ and b'_i is the number of infected machines on the next time tick. The recursion process will stop when there are no more vulnerable machines left. At some time tick, the total number of infected machines will be $b_1 + (2^8 - 1) \cdot b_2 + (2^{16} - 2^8) \cdot b_3$.

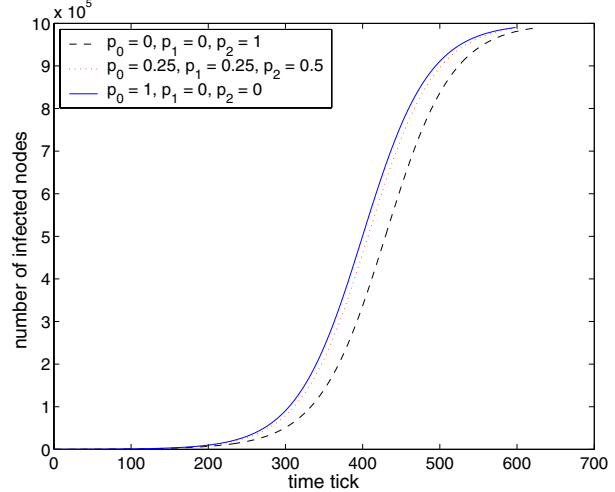
Based on the above formulae, we can understand the characteristics of local subnet scanning and the effect of the hitlist’s distribution. Different p_0 , p_1 , p_2 and h_1 , h_2 , h_3 can generate different patterns for the spread of worms.

Four cases are considered here:

- 1) Random scanning ($p_0 = 1$, $p_1 = 0$, $p_2 = 0$): In this case $k_1 = k_2 = k_3 = \frac{\text{number of total infected machines}}{2^{16}}$, which means the distribution of the hitlist cannot effect the spread of active worms.
- 2) A hitlist with an even distribution ($h_1 = h_2 = h_3$): This gives $k_1 = k_2 = k_3 = s b_1 = s b_2 = s b_3$. Local subnet scanning, therefore, cannot change the spread of active worms in this case.
- 3) Similar to the Nimda worm ($p_0 = 0.25$, $p_1 = 0.25$, $p_2 = 0.5$): In this case, we select different distributions of the hitlist, just as in Figure 7(a). Evenly distributed hitlists give the best performance, while putting all hitlists together in one “special” subnet ($h_1 = 10$, $h_2 = h_3 = 0$) gives us the worst performance. This figure shows that



(a) Local subnet scanning with different hitlist distributions. All cases are for $p_0 = 0.25$, $p_1 = 0.25$, $p_2 = 0.5$.



(b) Local subnet scanning with an uneven hitlist distribution. All cases are for $h_1 = 10$, $h_2 = \frac{40}{2^8-1}$, $h_3 = \frac{50}{2^{16}-2^8}$.

Fig. 7. Modeling the Spread of Active Worms that Employ Local Subnet Scanning (All cases are for 1,000,000 vulnerable machines which are evenly distributed to every subnet, a scanning rate of 100 scans/time tick and a time period of 1 time tick to complete infection.)

the hitlist's distribution can affect the spread of active worms.

- 4) Local subnet scanning with a hitlist of uneven distribution (fix h_1 , h_2 , h_3 and let $h_1 > h_2 > h_3$): This stands for a hitlist of uneven distribution and a centralization of more hitlist machines in the "special" subnet. Surprisingly, however, Figure 7(b) shows that in this case local subnet scanning slows down the propagation of active worms. We will further discuss why worm designers select this scanning technique in the next section.

From the four cases above, we see that for local subnet scanning the hitlist's distribution can influence the spread of active worms, while the even distribution gives us the best performance. In addition when the hitlist is more concentrated in the "special" subnet, local subnet scanning slows down the spread of active worms.

B. Discussion of the Local Subnet Scanning Policy

The LAAWP model implies that local subnet scanning may slow down the spread of active worms. Why do the designers of active worms use this technique? There are two main reasons:

- 1) Firewalls can protect vulnerable machines behind it. But local subnet scanning allows a single copy of a worm running behind the firewall to rapidly infect all the other local vulnerable machines.
- 2) One subnet always belongs to a company or organization and has a lot of similar machines. Therefore, it can be expected that if a machine has a security hole, then there is a high probability that many other machines in the same network have the same security hole.

VI. CONCLUSIONS

In this paper we present the AAWP model to analyze the characteristics of the spread of active worms. Even though the AAWP model also used deterministic approximation, it gives more realistic results when compared to the Epidemiological model. The simulation results show that our model can better explain the "mystery" in [9]. The AAWP model can be used to simulate the Code Red v2 worm with the following parameters: 500,000 vulnerable machines, starting on a single machine, a scanning rate of 2 scans/second, a death rate of 0.00002 /second, a patching rate of 0.000002 /second, and a time period of 1 second to complete infection.

Taking the Code Red v2 worm as an example, we apply our model to answer three different questions. First, from our model we assert that an address space of 2^{24} IP addresses is large enough to obtain realistic results, while an address space smaller than 2^{20} addresses is not large enough to effectively obtain any realistic information about the spread of worms. Second, the AAWP model is used to evaluate the performance of a simple sensor detection system. More than 2^{18} unused IP addresses are needed for the sensors to detect the Code Red v2 like worm in one hour. Worms with a larger hitlist can be detected in a shorter period of time, even though they spread at a much faster rate. This simple sensor detection system is the first step towards a practical detection system that detects active worms through scanning frequencies or source IP address distributions. We plan to use our model to evaluate this type of detection system. Finally, the AAWP model is used to evaluate the performance of the LaBrea tool defense system. Similarly, an address space of more than 2^{18} unused IP addresses is needed by LaBrea to defend against the Code Red v2 like worm

effectively. We plan to apply our model to assess other publicly available defense systems and compare the relative performance of different defense systems.

As part of our ongoing work, we extend the AAWP model to the LAAWP model to understand the spread of active worms using local subnet scanning. The distribution of the hitlist can affect the local subnet scanning policy. In particular, a worm using an evenly distributed hitlist spreads at the fastest rate. When the hitlist is concentrated in some subnet, the spread of active worms is slowed down. In the LAAWP model, the vulnerable machines are assumed to be evenly distributed in every subnet. We plan to study the effect of the distribution of vulnerable machines in order to get more accurate results.

REFERENCES

- [1] R. Russell and A. Machie, "Code Red II Worm," Tech. Rep., Incident Analysis, SecurityFocus, Aug. 2001.
- [2] A. Machie, J. Roculan, R. Russell, and M. V. Velzen, "Nimda Worm Analysis," Tech. Rep., Incident Analysis, SecurityFocus, Sept. 2001.
- [3] CERT/CC, "CERT Advisory CA-2001-26 Nimda Worm," <http://www.cert.org/advisories/CA-2001-26.html>, Sept. 2001.
- [4] D. Song, R. Malan, and R. Stone, "A Snapshot of Global Internet Worm Activity," Tech. Rep., Arbor Networks, Nov. 2001.
- [5] S. Staniford, V. Paxson, and N. Weaver, "How to Own the Internet in Your Spare Time," in *Proc. of the 11th USENIX Security Symposium (Security '02)*, 2002.
- [6] J. O. Kephart and S. R. White, "Measuring and modeling computer virus prevalence," in *Proc. of the 1993 IEEE Computer Society Symposium on Research in Security and Privacy*, May 1993, pp. 2–15.
- [7] J. O. Kephart, "How topology affects population dynamics," in C. Langton, ed., *Artificial Life III. Studies in the Sciences of Complexity*, 1994, pp. 447–463.
- [8] J. O. Kephart and S. R. White, "Directed-graph Epidemiological Models of Computer Viruses," in *Proc. of the 1991 IEEE Computer Society Symposium on Research in Security and Privacy*, May 1991, pp. 343–359.
- [9] S. R. White, "Open Problems in Computer Virus Research," presented at Virus Bulletin Conference, Oct. 1998.
- [10] N. Weaver, "Warhol Worms: The Potential for Very Fast Internet Plagues," <http://www.cs.berkeley.edu/~nweaver/warhol.html>.
- [11] D. Moore, "The Spread of the Code-Red Worm (CRv2)," http://www.caida.org/analysis/security/code-red/coderedv2_analysis.xml.
- [12] C. C. Zou, W. Gong, and D. Towsley, "Code Red Worm Propagation Modeling and Analysis," in *9th ACM Conference on Computer and Communication Security*, Nov 2002.
- [13] CAIDA, "CAIDA Analysis of Code-Red," <http://www.caida.org/analysis/security/code-red/>.
- [14] T. Liston, "Welcome To My Tarpit - The Tactical and Strategic Use of LaBrea," <http://www.hackbusters.net/LaBrea/LaBrea.txt>.
- [15] T. Liston, "LaBrea," <http://www.hackbusters.net/LaBrea/>.