

Capturing Ghosts: Predicting the Used IPv4 Space by Inferring Unobserved Addresses

Sebastian Zander
CAIA, Swinburne University of
Technology
Melbourne, Australia
szander@swin.edu.au

Lachlan L. H. Andrew
Faculty of IT,
Monash University
Melbourne, Australia
lachlan.andrew@monash.edu

Grenville Armitage
CAIA, Swinburne University of
Technology
Melbourne, Australia
garmitage@swin.edu.au

ABSTRACT

The pool of unused routable IPv4 prefixes is dwindling, with less than 4% remaining for allocation at the end of June 2014. Yet the adoption of IPv6 remains slow. We demonstrate a new capture-recapture technique for improved estimation of the size of “IPv4 reserves” (allocated yet unused IPv4 addresses or routable prefixes) from multiple incomplete data sources. A key contribution of our approach is the plausible estimation of both observed and unobserved-yet-active (ghost) IPv4 address space. This significantly improves our community’s understanding of IPv4 address space exhaustion and likely pressure for IPv6 adoption. Using “ping scans”, network traces and server logs we estimate that 6.3 million /24 subnets and 1.2 billion IPv4 addresses are currently in use (roughly 60% and 45% of the publicly routed space respectively). We also show how utilisation has changed over the last 2–3 years and provide an up-to-date estimate of potentially-usable remaining IPv4 space.

Categories and Subject Descriptors

C.2.3 [Computer-Communication Networks]: Network Operations—*Network Monitoring*; C.4 [Performance of Systems]: Measurement Techniques

General Terms

Measurement

Keywords

Used IPv4 space; capture-recapture

1. INTRODUCTION

At the end of June 2014 less than 4% of the IPv4 address space remained unallocated by Regional Internet Registrars (RIRs). RIPE and APNIC have exhausted their supply and the other RIRs (except AfriNIC) will run out of prefixes by the end of 2014 [1]. Understanding the pressures for IPv6 adoption, and the scope of possible

IPv4 address markets, requires plausible estimates of actual IPv4 address use – particularly the efficiency with which allocated prefixes are filled with actively-used addresses. Ideally, our estimation techniques should also help the community track progressive exhaustion once all routable IPv4 prefixes are allocated.

Prior studies that, among other things, analysed IPv4 space growth [2–4] and a port scan census from 2012 [5] used mainly active probing (“pinging”). Yet pinging alone will under-count, as many hosts do not respond or their responses are filtered (e.g., by firewalls). Recently, Dainotti *et al.* [6] used IPv4 data from multiple sources to estimate the used /24 networks. Apart from a simple multiplier in [3], previous work did not attempt to correct for under-sampling.

Our key contribution in this work is a new method to estimate the true population of both observed and unobserved (yet still active) IPv4 addresses using a statistical *capture-recapture* (CR) [7–9] model applied over diverse sources of active and passive measurement data. We significantly extend our earlier workshop paper, with refined methodology, additional data sources and greatly extended analysis [10].

Our second contribution is a three-year study of address use using our CR method. We “pinged” the allocated space with ICMP echo requests and TCP port 80 probes, and also gathered IPv4 data from web server logs [11], email spam detector logs [12], Wikipedia edit logs, logs of Valve’s Steam online game platform, logs from Measurement Lab [13], and university access router’s NetFlow logs. Inevitably, our sources only detect actively used addresses from 80% of the allocated space that is publicly routed (based on [14]). Hence, our analysis is focused on the routed space.

Although our sources provide diverse evidence of active IPv4 address use, there are likely many in-use addresses that we never see. We utilise our CR method to estimate a total population of used IPv4 addresses (and /24 networks) that *includes* these unobserved addresses (ghosts). As many sources obtain measurements over weeks or months, our estimates of the used IPv4 addresses (and /24 networks) are based on observation periods rather than points in time. By cross-validation with our datasets, and comparison with a few samples of ground truth, we show our CR method provides better estimates than prior techniques. We analyse “demand” – growth in address use – over the last 2–3 years relative to factors such as the RIR, country, or prefix size, and estimate the remaining “supply” of unused prefixes.

With just ICMP pinging we observed 4.9 million used /24 subnets and 430 million used IP addresses. Our combined sources observed 5.9 million used /24 subnets and 740 million used IPv4 addresses, yet our CR technique indicates significantly higher actual usage. We estimate 6.3 million /24 subnets and 1.2 billion IPv4 addresses were used by the end of June 2014 (approx. 60% and 45%

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than ACM must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from permissions@acm.org.

IMC’14, November 5–7, 2014, Vancouver, BC, Canada.

Copyright 2014 ACM 978-1-4503-3213-2/14/11 ...\$15.00.

<http://dx.doi.org/10.1145/2663716.2663718>.

of the publicly routed space respectively). From the end of 2011 to June 2014, the growth in used /24 subnets and IPv4 addresses was roughly linear, with an increase of 0.45 million /24 subnets and 170 million IPv4 addresses per year.

These trends mean routed but currently unused space could supply us until 2023. However, supply varies significantly across regions, for example Asia and South America only have supply for another 2–4 years (without reallocations between RIRs). Moreover, if, for example, only 75% of all routed /24 subnets could ever be used, regions with tight supply, like Asia and South America, have less than 1 year of supply left. Unrouted unused space may provide more supply.

The paper is organised as follows. Section 2 discusses related work. Section 3 describes the concept of CR and our log-linear CR models. Section 4 describes our IPv4 address data collection and processing. Section 5 covers the validation of our CR model. In Section 6 we analyse the growth of used IPv4 space over time, and in Section 7 we estimate the space still unused. Section 8 concludes and outlines future work.

2. RELATED WORK

The related measure of *routed* address space has been estimated based on prefixes advertised by BGP [15, 16]. However, estimation of the number of *actively used* addresses began with Pryadkin *et al.* [2], who used ICMP echo and TCP SYN probing to probe the allocated Internet. They discovered 62 million used IPv4 addresses in 2003 and 2004. Pryadkin *et al.* also showed that only a small number of allocated prefixes appeared to be heavily used, while a large part of the IPv4 space appeared unused or underutilized.

Heidemann *et al.* [3] infrequently probed all allocated IPv4 addresses (census) and frequently probed selected address samples (survey) with ICMP echo pinging to study usage, availability and up-time of addresses. The last census in 2007 from [3] accounted for 112 million used addresses. Heidemann *et al.* compared ICMP probing with TCP port 80 probing and passive measurements based on small samples. They proposed a correction factor of 1.86, thus estimating the total number of used IPv4 addresses in mid 2007 was 200–210 million.

Cai *et al.* [4] used ping survey data from [3] and conducted more surveys in 2009–2010 to analyse typical address block sizes and their characteristics. They did not directly estimate the used IPv4 address space, but observed: “most addresses in about one-fifth of /24 blocks are in use less than 10% of the time”.

From June to October 2012, anonymous researchers used hacked commodity routers to perform a port scan of the IPv4 Internet [5]. They detected 420 million addresses that responded to ICMP echo, which is broadly consistent with our two ping censuses that detected 360 million addresses between March and September 2012.

In 2013 we initially proposed using a log-linear CR model to estimate the true population of used IPv4 addresses from multiple sources of IPv4 addresses [10]. Our preliminary workshop paper found that our log-linear CR estimate is significantly higher (one billion used IPv4 addresses in mid 2013) than the aggregate number of observed IPv4 addresses from multiple measurement sources.

Dainotti *et al.* [6] developed techniques to filter out spoofed IPv4 addresses from darknet or NetFlow data and showed that the filtered datasets can be used to estimate Internet address space usage. With the filtered darknet data, NetFlow data, and ping census data from USC [3] combined, they estimated 4.8 million used /24 subnets (47% of the routed space) in September 2012. This is broadly consistent with the 5.2 million /24 subnets we observed in the year to September 2012 (c.f. Figure 4). The difference is likely due to the larger number of sources and the longer time window we use.

3. CAPTURE-RECAPTURE

There are many techniques for estimating population sizes from limited samples. Some use problem-specific approaches, but many use CR methods. CR methods have been used in ecology [7, 8], epidemiology [9, 17], and to estimate missing links from observed AS-graphs [18].

First, we discuss general assumptions for CR. Then, to illustrate CR, we discuss the simplest CR technique – the two-sample Lincoln-Petersen (L-P) method. Since we have more than two sources and for our data some assumptions of the L-P method are violated, we do not use this method. Finally, we describe the log-linear CR models that we use. Log-linear models make less restrictive assumptions and work with arbitrarily many sources.

3.1 General assumptions

A prime assumption of CR is that all individuals of the population can be uniquely identified. This assumption obviously holds for IPv4 addresses (we only care if an address was used but not who used it). Another assumption is that the data sources only sample “alive” individuals. We achieve that by filtering out IPs from the data sources that were sampled but not actually used, for example due to address spoofing (see Section 4).

Any individuals with zero sample probability are not part of the CR estimate. In our case these are all used IP addresses in publicly unrouted space, which our sources cannot sample. Hence, our CR estimates are only for the publicly routed space. Furthermore, there may be some specialised devices using public IP addresses, such as printers, that our current data sources also cannot sample. This means our results likely have a downward “bias”, but as discussed in Section 4.2 the error may be relatively small.

3.2 Two-sample Lincoln-Petersen method

The Lincoln-Petersen (L-P) method [7, 8] is ideal to illustrate the basic principle behind CR, but has restrictive assumptions that prevent us from using it.

3.2.1 Description

The two-sample L-P method works as follows. Given a first sample, that observes M individuals, the size of the population would be known if we knew what *fraction* of the population had been observed. To estimate this, L-P takes a second sample. Say it contains C individuals, of which R individuals occur in both samples. If the fraction of “recaptured” individuals in the second sample equals the fraction of the total population captured in the first sample, $R/C = M/N$, then the population N is [7, 8]:

$$N = \frac{MC}{R}.$$

In our context, the samples or “sources” are different active and passive measurements (see Section 4). For concreteness, consider Source 1 to result from pinging the entire IPv4 space and Source 2 to be all addresses in a server log. Based on the number of unique addresses observed by Source 1 and Source 2, and the number of unique addresses observed by both sources (Overlap) CR allows to estimate the number of unobserved addresses (Unseen), as illustrated in Figure 1.

3.2.2 Assumptions

The L-P estimate assumes that the probability of an individual being captured in one source does not depend on the probability of being captured in a different source (*independent sources*). It also assumes that, within a sample, each individual has an equal chance of being sampled (*homogenous population*). Furthermore, the L-P

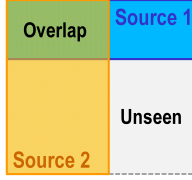


Figure 1: Two-source capture-recapture illustrated

estimate assumes that during measurement no individuals enter or leave the population (*closed population*). However, a violation of the last assumption is just another form of heterogeneity.

Given our data sources (see Section 4.1), there is no significant causal relationship to introduce source dependence. While some samples are dependent, i.e. IPs observed in traffic flows between the two NetFlow-monitored sites or between one NetFlow site and one of the logged sites (e.g. Wikipedia or the gaming site), their number is very small compared to the overall size of the datasets (less than 1%). However, the population is very heterogeneous; for example, servers are more likely to respond to ping, while client machines may be more likely to appear in certain traffic logs.

Nevertheless, heterogeneity gives rise to what is called *apparent source dependence*. For example, two sources that are biased towards client machines will *appear* to be positively dependent given another source that is less biased towards clients. In fact it is well known that heterogeneity and source dependence are often confounded and cannot be clearly separated [19]. Apparent source dependence must be treated similarly to source dependence.

If there is (apparent) dependence such that two sources are positively correlated, the L-P estimator underestimates the population size: $R/C > M/N$ and so $N > MC/R$. Similarly, if two sources are negatively correlated, the L-P estimator overestimates the population size. If the sign of the correlation is known, then L-P estimates can be used to identify plausible lower or upper bounds [17].

3.3 Log-linear Models

Just as L-P uses a second sample to estimate the fraction of the population of the first sample, a third sample can be used to compensate the correlation between the first two samples. This is the basic idea of log-linear models (LLMs) used for CR [17, 19, 20], which can model (apparent) source dependence among arbitrarily many sources.

3.3.1 Description

Let N be the unknown number of distinct individuals of the population. Let t denote the number of sources indexed by $1, 2, \dots, t$. For each individual, let s_1 to s_t be defined such that $s_i = 1$ if the individual occurs in sample i and $s_i = 0$ otherwise. Then the string $s_1 s_2 \dots s_t$ is called the “capture history” of the individual. The observed outcome of all measurements can then be represented by variables of the form z_s , which are the numbers of individuals with each capture history $s = s_1 s_2 \dots s_t$. These are assumed to be instances of random variables Z_s . Note that individuals with the capture history $00\dots 0$ are unobserved, and our goal is to estimate $Z_{00\dots 0}$. Table 1 illustrates the variables Z_s for each possible capture history for three sources ($t = 3$); “yes” means an individual was observed and “no” means an individual was not observed by a source.

For each history s , let $h(s)$ be the set of sources in which the individual occurs; for example, $h(101) = \{1, 3\}$. Define the indicator function $\mathbf{1}_A = 1$ if statement A is true and 0 otherwise. We can now write the following system of equations in 2^t variables

Table 1: Three-source contingency table showing all possible capture histories and number of unseen individuals Z_{000}

		Source 1			
		yes		no	
		Source 2		Source 2	
Source 3	yes	Z_{111}	Z_{101}	Z_{011}	Z_{001}
	no	Z_{110}	Z_{100}	Z_{010}	$Z_{000}=?$

$u, u_1, u_2, \dots, u_{12}, \dots, u_{23}, \dots$ up to $u_{12\dots t}$:

$$\log(\mathbb{E}(Z_s)) = \sum_{h \subseteq h(s)} u_h = \sum_h u_h \mathbf{1}_{h \subseteq h(s)}. \quad (1)$$

For example, for $t = 3$, the system is

$$\begin{aligned} \log(\mathbb{E}(Z_{ijk})) = & u + u_1 \mathbf{1}_{i=1} + u_2 \mathbf{1}_{j=1} + u_3 \mathbf{1}_{k=1} \\ & + u_{12} \mathbf{1}_{i=1 \wedge j=1} + u_{13} \mathbf{1}_{i=1 \wedge k=1} \\ & + u_{23} \mathbf{1}_{j=1 \wedge k=1} + u_{123} \mathbf{1}_{i=1 \wedge j=1 \wedge k=1}. \end{aligned}$$

The estimate of $Z_{00\dots 0}$ is then $\hat{Z}_{00\dots 0} = \exp(u)$. If we take $\mathbb{E}[Z_s] = z_s$ then this system has 2^t unknowns but only $2^t - 1$ equations, as $Z_{00\dots 0}$ is unknown. Hence it is customary to assume $u_{12\dots t} = 0$ [17]. As the number of sources t increases, this t -way dependency becomes decreasingly important.

The model with all u_h (the saturated model) is very sensitive to small values of Z_s . For example, a zero count for some capture history may give $\hat{Z}_{00\dots 0} = 0$, regardless of the other Z_s [17]. Furthermore, the larger t is, the higher is the random error for some Z_s . Including “noisy” parameters u_h in a model results in a poor predictive performance of the model (referred to as over-fitting).

Over-fitting is mitigated by “model selection” (see Section 3.3.2), in which some u_h are forced to 0, to reflect assumed independence between certain combinations of sources. For example, setting $u_{12} = 0$ indicates sources 1 and 2 are independent. With such incomplete models, the system of equations is overdetermined, and the maximum likelihood parameters u are typically used, based on the assumption that Z_s result from uniform random sampling and are hence Poisson distributed.

Assuming the Z_s are Poisson distributed is appropriate if the upper limit for the Z_s is unknown. However, we can bound Z_s by the size of the publicly routed IPv4 space. Hence, we use right-truncated Poisson distributions defined over $[0, l] \cap \mathbb{Z}$, where l is the upper limit. These improve estimates substantially for small strata, where the counters are relatively close to the limit (see Section 5.2), but otherwise make little difference.

3.3.2 Model selection

Model selection for an LLM consists of selecting which u_h will be assumed *a priori* to be 0. The goal is to select the least complex model with “adequate” fit of the observed (and by assumption) unobserved individuals [20].

A common approach is to minimize an “Information Criterion” (IC). Two common ICs are [21]:

$$\text{AIC} = 2k - 2 \log(L), \quad \text{BIC} = \log(M)k - 2 \log(L)$$

where L is the likelihood of the data given the assumed model, k is the number of free parameters of the model and M is the number of observed individuals. AIC is used more often, but each has merits [22]. Section 5 compares the BIC and the AIC for our data. We choose the simplest model m such that no other model n has $\text{IC}_n < \text{IC}_m - 7$ [21].

In our case, k is the number of non-zero u_h , but L is difficult to obtain. AIC and BIC assume that each source samples uniformly and so L is the likelihood of a Poisson model. If the number of samples is large, the central limit theorem indicates that substantial deviations from the mean have very low likelihood. In our case, as in [17, 19], the randomness comes largely from the choice of sources to monitor, which is hard to characterise but has substantially higher variance. Hence the Poisson assumption selects too complex a model.

We mitigate this overfitting using the simple heuristic of dividing all z_s by some integer d when calculating L . It remains to select d . If d is so large that any z_s gets rounded to zero, the LLM breaks down. The further heuristic of selecting d to be the largest number less than $\min_s z_s$ appears to work well (see Section 5.1).

3.3.3 Estimate range

Besides computing point estimates, we also compute estimate ranges (used in Section 5). We use the procedure in [23] to compute a $100(1 - \alpha)\%$ profile likelihood “confidence interval” (CI) for \hat{N} . Note that this is not a true confidence interval in our case, since it is based on the assumption that each sample is drawn randomly, resulting in a Poisson number of samples with each history. In contrast, our samples arise from different, not completely random sampling procedures. Hence we treat these “confidence intervals” as merely a useful heuristic indication of the sensitivity to modelling variations and we set $\alpha = 10^{-7}$ to obtain wide CIs.

3.3.4 Sampling zeros

Even with appropriate model selection, if the number of samples across all sources is low, we may have a large number of z_s that are near zero leading to unreliable estimates. In our case this only occurs for a few small countries or territories when stratifying by country code (see Section 3.4). Hence, we exclude country codes with fewer than 1000 IP addresses observed by all sources from the results in Section 6.2 (where they are negligible) and in Section 6.7.

3.4 Stratification

We obtain more insight and also initially hoped to mitigate heterogeneity by stratifying the population in different ways. We classified IPv4 addresses as statically or dynamically assigned using the approach described in [10], and based on allocation and whois data we stratified by RIR (e.g. APNIC), country, prefix size, industry¹ and allocation age.

4. DATASETS AND PREPROCESSING

An IPv4 address is considered *used* if it responds to active probes or participates in connections. A *used* /24 subset contains one or more used addresses. This section describes our sources of used IPv4 address data, our data collection and processing, and our handling of both spoofed and dynamically assigned addresses.

4.1 Datasets

Our first two datasets are from actively probing the whole allocated IPv4 Internet using ICMP echo requests (**IPING**) and TCP SYN packets to port 80² (**TPING**). Since mid-2011 we probed each allocated IPv4 address (a census) once every 6 months. The first two censuses used ICMP probing and the rest used both ICMP and TCP probing (with TCP probing seeing over 7% more observed

¹“Industry” indicates whether address space is education, military, government, corporate, or ISP. We classified 88% of the allocated address space based on whois information (down to /17 networks).

²Initially we probed a sample of the Internet using different commonly used TCP ports and found port 80 to be the most responsive.

IP addresses). We limited the overall ping rate and used reversed bit counting for “traversing” the IP space. On average our prober sent only one packet every two hours to individual /24 networks, to minimise congestion, stay below typical ICMP or TCP rate limit thresholds and avoid triggering monitoring systems (on average we received only 10–20 complaints per census). For the first half of 2011 we use ICMP ping data collected by USC/LANDER [24].

Passively observed IPv4 data includes addresses from Wikipedia’s page edit histories³ (**WIKI**), potential spam email senders from [12] (**SPAM**), addresses of clients tested by Measurement Lab [13] tools (**MLAB**), web clients participating in our IPv6 readiness test [11] (**WEB**), anonymized server logs of game clients connecting to Valve’s Steam online gaming platform (**GAME**), and NetFlow records from *incoming* traffic of Swinburne University of Technology’s access router (**SWIN**)⁴ and Caltech’s access router (**CALT**).

We utilise data gathered from 2011 onwards. We generate datasets of unique /24 subnets by processing the IPv4 datasets and setting the last octet of each address to zero and then filtering out the duplicates. Table 2 shows the number of unique IPv4 addresses and /24 subnets per dataset for the years 2011–2013 (IPs for GAME omitted for confidentiality). Note that the numbers in the table cannot be used as growth trends due to sample method variations.

4.2 Host types sampled

We collected data from diverse locations, but CR estimates will only be useful if (1) our datasets sample all types of devices using public IPv4 addresses and (2) a type of device can be sampled by multiple datasets. We now discuss whether this is the case based on grouping devices into routers, servers/proxies, clients (e.g. PCs, smart phones), and specialised devices (e.g. printers, cameras).

ISP routers are sampled by IPING and TPING and may also appear in SWIN and CALT. Home routers are sampled by IPING and TPING (we confirmed that some responses came from Cable/DSL routers by inspecting web pages from IPs responding to TPING) and by all other sources (with NAT packets sent from home networks appear to come from home routers). Servers/proxies are sampled by IPING, TPING, SWIN and CALT. They can also appear in WIKI, SPAM and WEB. Clients are sampled mainly by WIKI, SPAM, MLAB, WEB, GAME, SWIN and CALT, but also appear in IPING. NAT’ed clients also appear in IPING and TPING. Specialised devices may be sampled by IPING and TPING.

Overall, our datasets sample most groups well, especially servers and clients, which we assume are the largest groups. Specialised devices are likely severely under-represented, but these are very hard to sample. The authors of [5] probed the entire IPv4 Internet on several hundred ports and detected 36 million addresses that only responded to TCP SYNs but not to ICMP. In our censuses 15–20 million IPs reacted to port 80 TCP SYNs but not to ICMP. The difference of 15–20 million addresses could be specialised devices listening only on specific ports⁵ we missed, but this number is small compared to our total estimate of 1.2 billion used addresses.

4.3 Time windows

We collected data from 1 Jan 2011 until 30 June 2014. To analyse the growth trend of used IPv4 addresses, we split our data into *overlapping* 12-month windows. Windows start every three months, so the first window starts at 1 Jan 2011 and the last window starts at 1 Jul 2013. The last window ends at 30 June 2014. This is a suitable trade-off between temporal resolution and noisy estimates.

³Modification time and IPv4 address of edits by unregistered users.

⁴Excluding all traffic flows of our active prober.

⁵E.g., printers responding only on the Internet Printing (IPP) port.

Table 2: Data sources and observed unique IPv4 addresses and /24 subnets per year (SWIN and CALT after spoofed IP filtering)

Dataset	Description	Time collected	2011		2012		2013	
			IPs [M]	/24 [M]	IPs [M]	/24 [M]	IPs [M]	/24 [M]
WIKI	Wikipedia’s page edit histories	Jan 2011 – Jun 2014	5.5	1.69	5.9	1.97	6.8	2.16
SPAM	Potential spam email senders	May 2012 – Jun 2014	-	-	19.2	1.56	17.5	1.73
MLAB	Clients tested by Measurement Lab	Jan 2011 – Jun 2014	30.0	2.66	27.6	2.69	21.5	2.49
WEB	Web clients tested for IPv6	Mar 2011 – Jun 2014	22.0	2.92	88.0	3.89	108.7	4.13
GAME	Game clients logged into Valve’s Steam	Jan 2011 – Jun 2014	conf	3.11	conf	3.62	conf	4.32
SWIN	Swinburne access router NetFlow records	Jan 2011 – Jun 2014	150.6	3.13	142.4	3.38	112.9	3.36
CALT	Caltech access router NetFlow records	Jun 2013 – Jun 2014	-	-	-	-	356.8	3.92
IPING	ICMP ping census of IPv4 Internet	Mar 2011 – Jun 2014	320.3	4.24	358.2	4.54	411.1	4.81
TPING	TCP port 80 census of IPv4 Internet	Mar 2012 – Jun 2014	-	-	70.0	3.38	92.7	3.71

Also, for some datasets we cannot have smaller gaps between windows, e.g. we only conducted IPING/TPING censuses every six months and we only collected GAME data every 3+ months.

Overlapping windows smooth out quick changes, but we believe fast transients in the number of used IPv4 addresses are unlikely. In the rest of the paper we associate statistics with the end of time windows. For example, for the first window the observed and estimated used space is associated with 31 December, 2011.

4.4 Data collection and processing

Our goal is to get datasets of publicly routed IPv4 addresses that were actually used. We do not distinguish whether IPs were used legitimately or illegitimately, such as addresses hijacked by spammers. For active probing we assume that positive responses only occur when IP addresses were used – the same assumption made by previous work [2–4].

For IPING we only counted IPv4 addresses that returned ICMP echo replies, “destination protocol unreachable” or “destination port unreachable” messages. We ignored addresses with other ICMP errors or “TTL exceeded” messages, as for these it is unclear if they were actually used. For TPING we only counted addresses that returned SYN/ACKs. We ignored addresses that returned RSTs, as 25% of RSTs cover nearly contiguous /25 or larger networks, suggesting they may have originated from firewalls. Lack of reply indicates an address was truly unused, a host ignored the probe, or the probe or response was filtered or lost.

For the passive datasets we extracted the IPv4 addresses from log files. We filtered out multicast and private addresses (e.g., 10.0.0.0/8), and those in unallocated or unrouted space. We identified the routed space based on the route-views dataset from RouteViews (RVs) [14]. For each time window we downloaded weekly snapshots from RV and then aggregated all the snapshots (excluding a few unallocated but advertised prefixes).

For WIKI, SPAM, MLAB, WEB and GAME (server logs), the addresses are only recorded for successful TCP sessions, so we can be sure the addresses were used. SWIN and CALT contain unused IPs from inside the two networks due to scanning that we cannot filter out, but their number is negligible compared to the sizes of the datasets. More problematically is that SWIN and CALT also contain spoofed IPv4 addresses from outside the two networks that do *not* represent used addresses. We describe our filtering of these in the following section.

4.5 Removal of spoofed IPs

We only have lists of observed IP addresses for SWIN and CALT. We do not have packet or flow data, and we do not know from

which IPs flows originated. Hence, instead of using the technique in [6], we needed a new heuristic to remove spoofed addresses.

Our heuristic is based on the assumption that spoofed IP addresses, which were not actively used, are uniformly distributed over the IP space, since there are two main reasons for these:

- Distributed Denial of Service (DDoS) attacks where the attacking machines send traffic from spoofed source addresses selected randomly [25].
- Port scans where the scanner uses decoys. For example, the prominent nmap scanning tool has an option that allows creating additional scans from random spoofed addresses to obfuscate the identity of the scanning host(s).

Note that for CR spoofed addresses only cause errors if *they were not actually used*. DDoS attacks where the attacker spoofs the source address to be the victim’s address (reflector attacks) are not a problem, as for these the spoofed addresses are most likely used. Also note that while the autonomous systems (ASs) filtering spoofed source addresses may not be uniformly distributed, this does not affect the overall uniform distribution, since attackers or scanners (or their machines) are usually distributed across the whole Internet *including* ASs that do not filter.

The assumption of uniform distribution is supported by circumstantial evidence from our data. We observed that the unfiltered SWIN and CALT have uniformly random distributed IPv4 addresses in six /8 prefixes that were completely or almost completely unused by other sources (e.g. 53.0.0.0/8 or 55.0.0.0/8).⁶ While the number of observed IPs from these ‘empty’ /8 subnets differs for SWIN and CALT, for a given dataset and time period the number of observed IPs is roughly identical for these /8 – consistent with the assumption that spoofed addresses are uniformly distributed over the IPv4 space.

Our approach works in two stages. First, we estimate which /24 subnets should be removed entirely, and then we remove potentially spoof addresses from used /24s.

From SWIN and CALT we removed all /24 subnets that:

1. have fewer than m observed IPs, *and*
2. have no overlapping IPs that are also in the spoof-free WIKI, WEB, MLAB and GAME datasets.

⁶The number of addresses from these /8 in our non-spoofed sources is negligible (no more than a few tens of addresses) and in some cases we know from the network administrators that these /8 are hardly used. However, for SWIN and CALT we see more than 10,000 addresses in these /8.

We choose m as follows. Treating spoofed IPs as uniformly sampled from a space of s IPs with probability p , the number X of spoofed IPs in the space follows a Binomial distribution. Specifically:

$$\Pr(X > k) = 1 - \sum_{i=0}^k \binom{n}{i} p^i (1-p)^{n-i}.$$

In our case of $/24$ subnets, $s = 256$ and we estimate p based on the number of spoofed IPs S in each ‘empty’ $/8$ prefix, so $p = S/2^{24}$. We then choose $m = k$ where $\Pr(X > k) < 10^{-8}$. Note that for SWIN, S is relatively constant across all time periods (10,000–15,000), but for CALT it increases from 15,000–20,000 until December 2013 to almost 250,000 in March 2014.

Spoofed IP addresses will also fall into $/24$ subnets that have actually used IP addresses. The second phase is to filter out potentially spoofed IPs in used $/24$ as follows. Since we assume the spoofed IPs to be uniformly random distributed, the number of spoofed IPs is S for used $/8$ prefixes as well. Subtracting the number of already removed IPs in spoofed $/24$ subnets we have S'_i spoofed IPs left in $/8$ prefix i . Given the observed number of IPs T_i in $/8$ prefix i in SWIN or CALT the expected number of not-spoofed addresses per $/8$ prefix (out of 2^{24} addresses) is

$$2^{24} \cdot \frac{T_i - S'_i}{2^{24} - S'_i}.$$

On average the probability that an IP in i is valid (V) is

$$\Pr(V) = \left(\frac{T_i - S'_i}{T_i} \right) \left(\frac{2^{24}}{2^{24} - S'_i} \right) \approx (T_i - S'_i) / T_i.$$

This tells us how many IPs to keep in each $/8$ prefix, but we must also determine which IPs to keep. To do this we use the fact that the distribution of the final byte B of used addresses is not uniform. We estimate the probability $P(B|V)$ from the IPs observed by all sources except SWIN and CALT. Then assuming that $P(B|\text{not } V) = \frac{1}{256}$ (uniform distribution), Bayes’ rule gives that an IP is not spoofed in SWIN or CALT with probability

$$\Pr(V|B) = \frac{\Pr(V) \Pr(B|V)}{\Pr(V) \Pr(B|V) + (1 - \Pr(V)) / 256}.$$

We then filter SWIN and CALT by independently removing addresses ending with B with probability $1 - P(V|B)$.

We cannot evaluate the true accuracy of our approach, but the following circumstantial evidence shows that it is effective. With filtering, randomly distributed IPs in the ‘empty’ $/8$ networks are removed. With filtering, the number of used $/24$ subnets gradually increases over time and does not show large abrupt increases and decreases anymore. Without filtering the number of $/24$ subnets in SWIN or CALT is much higher than in any other dataset, e.g. it is up to 30% higher than for our largest dataset (IPING) and up to 60% higher compared to WEB, GAME. After filtering the number of used $/24$ subnets in SWIN and CALT is lower or similar to that in WEB and GAME.

Figure 2 shows the benefit of filtering spoofed addresses. LLM estimates that include filtered SWIN and CALT are quite consistent with LLM estimates made without SWIN and CALT. LLM estimates using unfiltered SWIN and CALT are much higher (exceeding the possible maximum for March 2014). To save space, we only show this comparison for $/24$ subnets, as spoofed IPs have less negative impact on the number of observed and estimated used IPv4 addresses due to the uniform random nature and low (10% or less) estimated percentage of spoofed IPs.

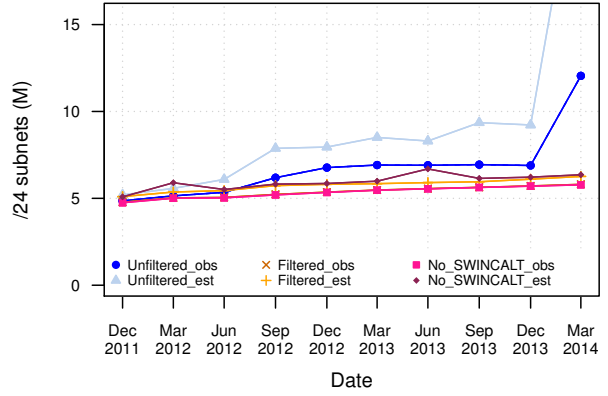


Figure 2: Observed (obs) and estimated (est) $/24$ subnets with and without spoof filtering compared to observed and estimated $/24$ subnets without SWIN and CALT

4.6 Dynamic and static addresses

Many IPv4 addresses are (re)assigned dynamically (such as with DHCP or PPPoE). Hence, long passive measurements may observe multiple addresses for a single host, and over-count the number of simultaneously used addresses.

If each assignment uses the lowest/highest unused address of a pool, then the total number of addresses used from the pool is the maximum simultaneous pool utilisation and the LLM estimate would indeed estimate the maximum number of simultaneously used addresses. However, if addresses are drawn uniformly, as our measurements suggest, then all pool addresses could be observed even if at most one address is in use at a time. Similarly, a single host moving between multiple statically assigned addresses may report multiple addresses, even if at most one is in use at a time.

However, addresses assigned to pools cannot be used elsewhere. So we argue that any over-count captures addresses (or $/24$ subnets) that are on ‘stand-by’ and de facto ‘in use’ at the time of our measurement. (In the future under-utilised pools may be reduced in size and the freed addresses may be used for other purposes. However, this is the same as re-purposing addresses of de facto unused hosts. We cannot quantify such future optimisations.)

We also study $/24$ subnets, which are less affected by dynamic addressing [6]. While address reassignments (e.g., host mobility) may cross different $/24$ subnets, a large fraction of them will be within the same $/24$ subnets. Evidence for this we found in 16 consecutive days of GAME session data, where we have unique client IDs, client login/logout times and client IP addresses. We selected 9 million distinct clients with multiple sessions and analysed the number of distinct IPv4 addresses and $/24$ networks used over time. After the first four days all clients had logged in at least once. From this point in time the observed distinct IP addresses increased 2.7 times (from 16 to 42 million), while the observed distinct $/24$ networks only increased 1.2 times (from 2.3 to 2.8 million).

5. VALIDATION

In this section, we first pick a specific model-selection algorithm from among those described in Section 3.3.2, based on test data. Then, we compare estimated use of addresses against ground truth for a handful of networks, and show that CR gives better estimates than simply summing the observed addresses. Finally, we use *cross-validation* to demonstrate that this also applies to the whole address space, since we have no ground truth for most networks.

Table 3: Cross-validation errors depending on different model selection parameter settings

Setting	IP addresses		/24 subnets	
	RMSE [M]	MAE [M]	RMSE [k]	MAE [k]
AIC-fixed1	28.1	11.9	139.1	57.2
BIC-fixed1	31.7	13.3	136.6	56.4
AIC-fixed10	20.6	9.5	138.4	57.7
AIC-fixed100	11.9	6.6	141.9	59.8
AIC-fixed1000	18.9	9.8	132.5	61.4
AIC-adaptive1000	15.5	7.7	134.9	56.3
BIC-adaptive1000	14.9	7.8	136.6	56.9

To perform cross-validation with our $k = 9$ data sources we consider a particular source i as the “universe” of possible IPv4 addresses. We apply CR to the addresses/subnets in i that are also in the other $k - 1$ sources, to estimate the number of individuals unique to source i . Since we know the true number of individuals unique to i , we can evaluate the effectiveness of CR. We do this for each source, to obtain the mean error and mean-square error. We then assume that the CR estimator based on the full k sets is equally accurate at estimating the true number of ghosts, although we do not have confidence intervals for this.

5.1 Model selection

We first investigated how to best select models, as our model selection approach in Section 3.3.2 leaves us with the choice of two ICs and the choice of how to dimension the count pre-processing heuristic. In initial experiments we varied the IC used (AIC, BIC) and we varied d for the count pre-processing heuristic. The variable d was either *adaptive* (starting with $d = 1000$ in each step we halved d , selecting the first d smaller than $\min_s z_s$) or *fixed* to values of 1, 10, 100, or 1000.

For each parameter setting we performed cross-validation for each time window, except the first window for which we have the least data, for both used IPv4 addresses and used /24 subnets. Table 3 shows the different parameter settings we investigated and the Root Mean Square Error (RMSE) and Mean Absolute Error (MAE) averaged over all sources and time windows.

With a fixed divisor using the actual counters (fixed1) results in the highest errors for IPs but provides low errors for /24 subnets. We think this is because (1) there is much more noise in the observed IPs than in the observed /24 subnets and (2) the number of observed /24 subnets is much lower than the number of IPs, and hence even for a small divisor of 10 we already start losing information for /24 subnets, which leads to much reduced accuracy. Choosing a divisor around 100 results in the smallest error for IPs but the largest error for /24 subnets. Effectively, the choice depends on the type of data and it is unclear what choice would be the best for estimating the IP addresses and /24 subnets unseen overall.

In contrast, our adaptive approach (with a maximum divisor of 1000) works quite well for both IPs and /24 subnets, with errors not much larger than the minimum errors. With the adaptive approach, using the BIC instead of the AIC lowers the error for IPs and increases the error for /24 subnets, but the increase for /24 subnets is small and even for /24 subnets the estimates obtained with the BIC are smoother.⁷ Hence, in the rest of the paper the estimates pre-

⁷The BIC selects fewer parameters representing interactions of many sources, which have a much lower number of samples than interactions between fewer sources and hence are noisier.

Table 4: Pingable, observed (obs), and estimated (Poisson, right-truncated Poisson) IPv4 addresses vs. peak usage (truth) as percentages of the size of each routed network

Network	Ping [%]	Obs. [%]	Estimated(Error) [%]		Truth [%]
			Poisson	TruncPoisson	
A	0.4	5.7	23.2(-2.7)	26.7(+0.8)	25.9
B	6.7	8.5	13.3(+1.9)	12.3(+0.9)	11.4
C	12.0	13.7	37.6(-)	36.1(-)	30–35
D	24.0	31.8	41.8(-5.8)	51.6(+4.0)	47.6
E	9.4	17.3	52.1(-6.2)	60.5(+2.2)	58.3
F	0.0	15.9	20.2(-2.1)	20.2(-2.1)	22.3

sented are based on our adaptive approach with a maximum divisor of 1000 and we use the BIC.

5.2 Comparison with ground truth

We compared our observed and estimated IP addresses with the ground truth for several networks. Our ground truth is estimates of the number of actively used IPv4 addresses at peak times (effectively high watermarks). Since our time windows are very long (12 months), it is appropriate to compare our estimates with high watermarks. We compare the ground truth with the observed and estimated numbers for the time windows where the high watermarks occurred roughly in the middle of the windows.

For privacy reasons we cannot reveal the identity of the networks. We also cannot reveal their sizes, as this would leak information, which allows narrowing down the possible identities (to a very small set in the case of one particular network). The largest network covered is two /16 subnets and the smallest network is roughly one /20 combined from multiple allocations.

For each network Table 4 shows the number of addresses that responded to ping, the number of addresses observed, the number of addresses estimated (for both Poisson and right-truncated Poisson), and the actual number of used addresses as percentages of the sizes of the routed networks. Note that network F blocked our pinger, so we do not have IPING or TPING data for this network.

The results show that the percentage of pingable and observed addresses is much smaller than the ground truth for networks A, C, D, E and F, whereas for the more “open” network B the percentage of observed addresses is relatively close to the truth. However, the CR estimates are always much closer to the truth. Using right-truncated Poisson distributions gives better estimates than using Poisson distributions. The right-truncated Poisson estimates tend to be higher than the truth (except for network F where we do not have IPING and TPING), but the cause may be dynamic addresses (e.g. DHCP) leading to higher estimates due to our long 12-month observation windows in contrast to the short-term peak numbers we use as truth.

5.3 Cross-validation results illustrated

Figure 3 illustrates the results of the cross-validation for addresses and subnets for time window 9 (results for other time windows are largely consistent). The figure shows the number of IPs in each source also observed by IPING (Obs ping), the total number of addresses of a source also observed by any other sources (Obs. all), and the ranges of the CR estimates (confidence interval based on profile likelihood with $\alpha = 10^{-7}$ to get wide intervals). Since the sources are of different sizes we normalised the number of addresses on the total number of addresses observed by each source

Table 5: Observed and estimated used IPv4 addresses and /24 subnets at the end of June 2014 based on different stratifications

Stratification	Estimated total [M]							Ping [M]	Observed [M]	Est. unseen [M]	Routed [M]
	None	RIR	Country	Age	Prefix size	Industry	Stat/Dyn				
IP addresses	1133.3	1167.2	1157.1	1107.0	1091.5	1084.8	1100.4	428.9	739.2	345–455	2725.3
/24 subnets	6.2	6.3	6.3	6.2	6.2	6.3	6.2	4.9	5.9	0.3–0.4	10.7

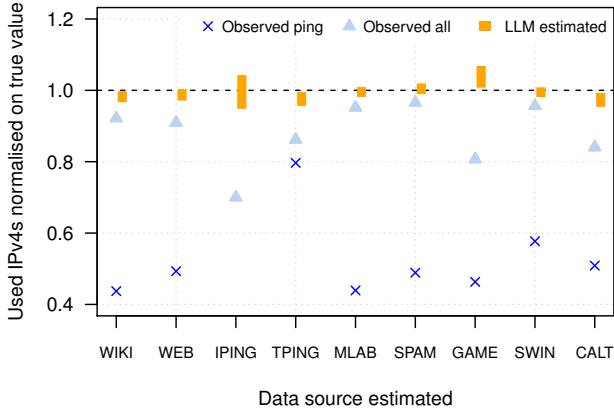


Figure 3: IP addresses observed with ping, observed by any source, and estimated ranges for LLM normalised on the true number of unseen IPs for each data source

(the ground truth). A CR estimate is good, if the normalised range includes 1 and the range is not too large.

Figure 3 shows that for IPv4 addresses all sources other than IPING and GAME have relatively high overlap, but between 10% and 15% of addresses appear only in one source. Only 50–60% of addresses of each source (except TPING) is in IPING, showing that ICMP ping undercounts significantly. The CR estimates for WIKI, WEB, IPING, MLAB, SPAM, and SWIN are quite good. The estimated ranges for TPING and CALT are slightly too low, and the estimated range for GAME is slightly too high. Nevertheless, the LLM CR estimates are a substantial improvement over just using the number of observed IPs.

For brevity we do not show the figure for /24 subnets. For /24 subnets there is a very high overlap between all data sources. However, for most sources only 90% or less of the /24 subnets appear in IPING, so just using ICMP ping significantly undercounts even the used /24 subnets. While the difference between CR estimates and observed addresses is much smaller for /24 subnets (in most cases the difference is only 1–2%, except for IPING), our CR estimates are still an improvement.

6. USED SPACE ANALYSIS

Now, we present the results for the estimated used IPv4 addresses and /24 subnets. We present both total estimates as well as estimates for different RIRs, countries, allocation ages, and allocation prefix sizes. We also investigate growth on a longer time scale and whether our estimates are sensible given the reported growth of Internet users.

6.1 Definitions

In most figures we plot one or more of the following metrics:

- *Routed IPv4 addresses and /24 networks*: the total number of addresses/networks that were publicly routed based on

RV [14]. For each time window we downloaded and aggregated weekly snapshots (excluding unallocated but advertised prefixes). Gregori *et al.* [26] suggests that while RV data is incomplete for AS-level graph analysis, it does capture the whole routed space.

- *Observed IPv4 addresses and /24 subnets*: the addresses/networks present in one or several of the data sources described in Section 4.1.
- *Estimated IPv4 addresses and /24 subnets*: the total number of addresses/networks estimated with CR, which means the observed addresses/networks plus the unseen estimated addresses/networks. For clarity in the figures we plot the point estimates instead of the estimated ranges.

In the figures with normalised data, we always normalise each series on the first value – the first time window (31 December 2011).

6.2 Used IPv4 space totals

Table 5 shows the estimated used IPv4 addresses and /24 subnets depending on different stratifications (RIR, country, prefix size, industry, allocation age) as introduced in Section 3.4, as well as the pingable, observed and estimated unseen addresses and /24 subnets at the end of June 2014. For each type of stratification we separated each source into the different strata, then used CR to estimate the size of each strata, and finally we summed up the estimates over all strata to get the total estimate.

Only 430 million IPs and 4.9 million /24s responded to ICMP ping, but we observed 740 million IPs and 5.9 million /24s from all sources combined. The estimated used IPs are fairly consistent across stratifications: roughly 1.1–1.2 billion used IPv4 addresses and 6.2–6.3 million used /24 subnets. Based on RV this means we observed only 27% of the routed IPv4 addresses and 55% of the routed /24 subnets, but we estimate that roughly 45% of the routed IPv4 addresses and 60% of the routed /24 subnets were used.

For all stratifications our estimates are always plausible (below the number of routed addresses). The quotient of estimated used addresses divided by the addresses detected only with ICMP echo ping is 2.6–2.7, larger than the correction factor of 1.86 in [3].

6.3 Used IPv4 space over time

Figure 4 shows the number of estimated used /24 subnets against the number of observed and routed /24 subnets both as absolute numbers and normalised. The dashed line is the actual estimates and the solid line is the estimates smoothed. The total number of observed /24 increased from 4.8 million to 5.9 million, but we estimate that the number of /24 subnets actually increased from 5.1 million to over 6.2 million (an increase of 0.45 million subnets per year). The estimate range is narrow, the minimum and maximum are within $\pm 1\%$ of the point estimates. Whilst the routed space only increased by 7% in two years, the number of observed and estimated used /24 subnets increased by 22% over the same time.

Figure 5 shows the number of estimated used IPv4 addresses against the number of observed and routed IPv4 addresses both as

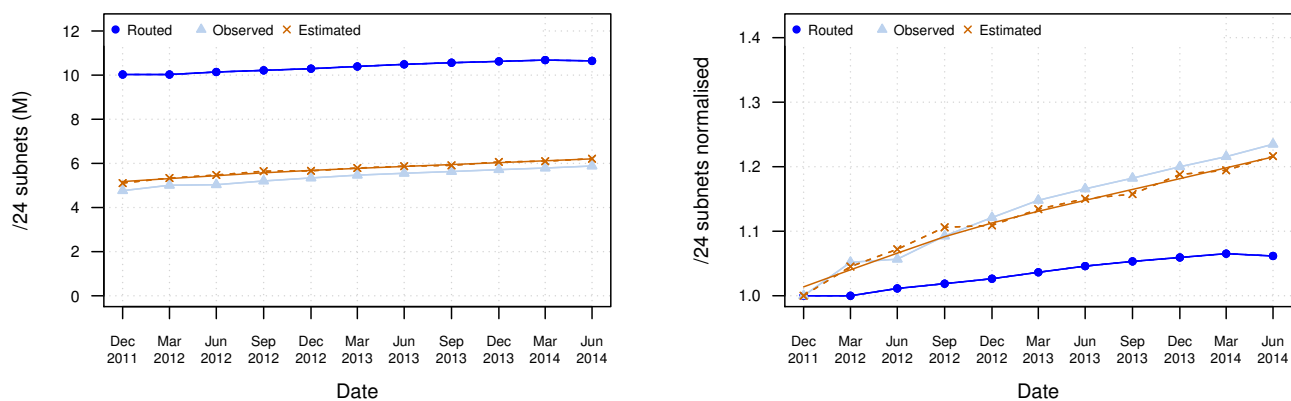


Figure 4: Absolute and relative growth of estimated, observed and routed /24 subnets

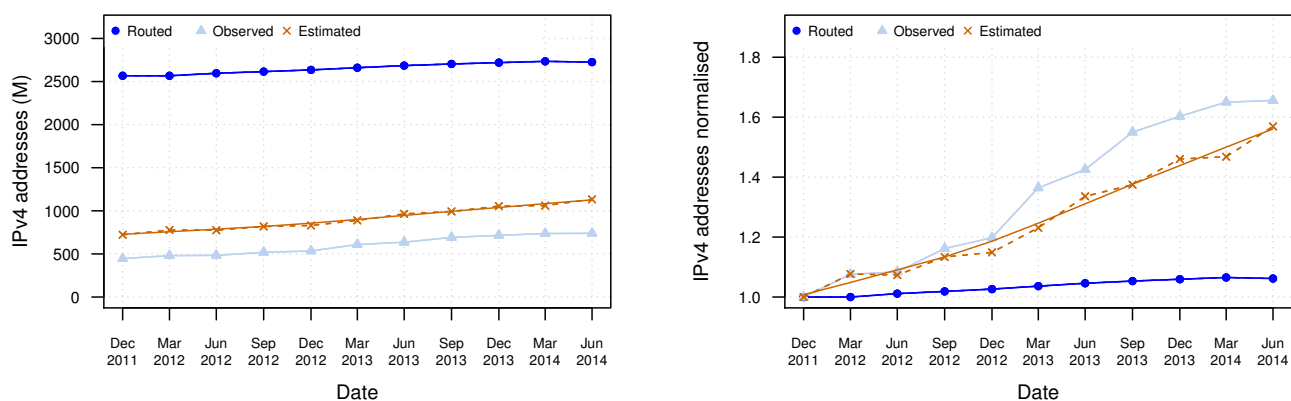


Figure 5: Absolute and relative growth of estimated, observed and routed IPv4 addresses

absolute numbers and normalised. The number of observed IPv4 addresses increased from 450 million to 740 million, but we estimate that the number of addresses actually increased from 720 million to 1.2 billion (an average increase of about 170 million IPv4 addresses per year). Minimum and maximum of the estimate range are within $\pm 3\%$ of the point estimate. As for /24 subnets, the observed and estimated number of IPv4 addresses increased faster than the routed addresses. The difference between estimated and observed relative growth may be in part because of earlier undercounting due to fewer sources and a gap in the GAME data collection.

The number of estimated /24 networks is only 5–10% above the number of observed /24 networks, whereas the number of estimated IPs is 50–60% above the number of observed IPs. Intuitively this makes sense, since a /24 network is observed if *any* of its addresses is observed. The relatively small error for /24s means one could treat the number of observed /24s as relatively good approximation of the number of actually used /24s.

6.4 Used IPv4 space by RIR

Figure 6 shows the estimated number of IPv4 addresses over time depending on the RIR responsible for their allocation both as absolute numbers and normalised. For brevity we omitted the broadly similar statistics for /24 subnets. APNIC has the largest number of used addresses followed by RIPE and ARIN. Looking at relative growth, AfriNIC is growing at the fastest rate, followed

by LACNIC. Of the three RIRs with the most allocated space, relatively APNIC and ARIN are growing faster than RIPE.

6.5 Used IPv4 space by prefix size

Figure 7 shows the average yearly growth rate for addresses for different prefix sizes (based on the RIR allocation data). For brevity we do not show the estimates for /24 networks here, as the trends are broadly similar. Absolute growth is strongest in the large prefixes /10 to /16 (/8 and /9 have not grown much). However, if we look at relative growth, growth has been more equally across many prefixes. Exceptions are the old /8 allocations which have not grown and /9, /21 and /22 allocations which show the strongest growth (/9 is driven up by a few ISPs since there are less than ten /9 allocations overall, and /22 is the largest allocation handed out by APNIC since 15 April, 2011, and by RIPE since 14 Sep, 2012).

6.6 Used IPv4 space by allocation age

Figure 8 shows the average yearly growth rate of IPv4 addresses for different allocation ages (based on RIR allocation data) until the end of 2013. For brevity we omitted the results for /24 subnets as the trend is broadly similar. In absolute numbers the more recent allocations made since 2005 are growing the most, with a clear positive correlation between recency and growth. In relative terms growth is strongest for allocations made in the last three years, but we can also see 20% or higher growth in some old allocations.

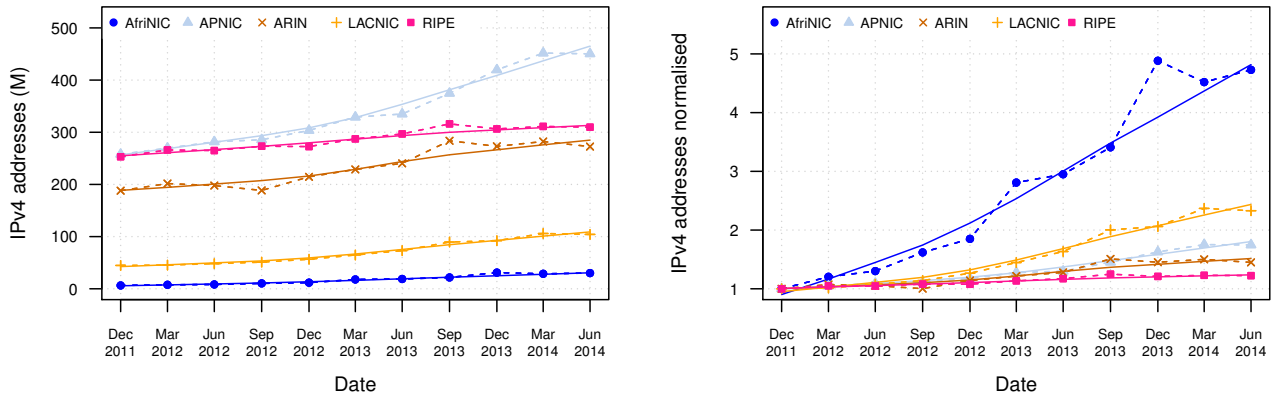


Figure 6: Absolute and relative growth of estimated IPv4 addresses for different RIRs

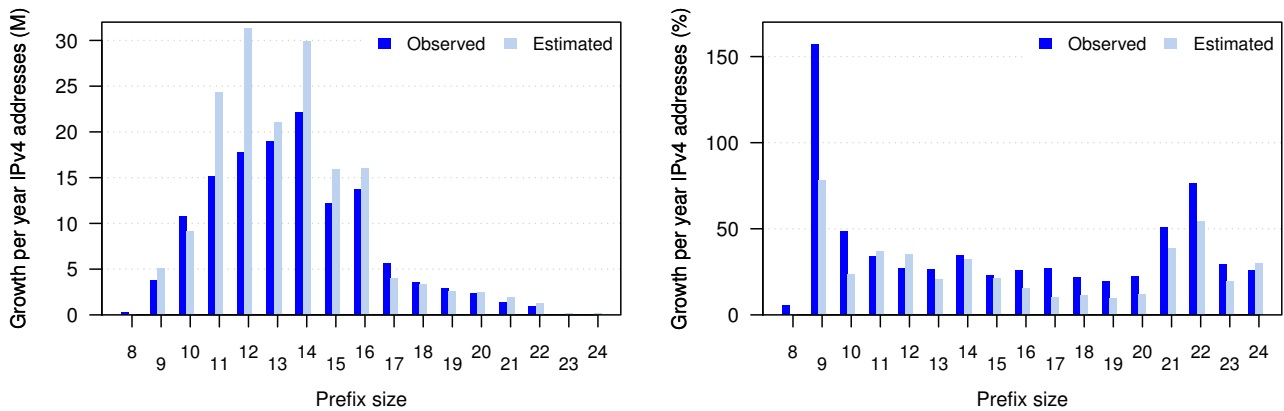


Figure 7: Average absolute and relative yearly growth of observed and estimated IPv4 addresses for different allocation prefixes

6.7 Used IPv4 space by country

Figure 9 shows the absolute and relative growth for IPv4 addresses for the countries with the largest number of observed used IPv4 addresses (at least 1.5 million addresses). Again, we do not show results for /24 subnets as the trends are broadly similar. Absolute growth is strongest in the two nations with the largest allocations (USA, China) followed by Brazil and South Korea. Relative growth is between 10% and 30% for many countries, but Romania and several Asian and South American countries (Brazil, Columbia, Indonesia, India, Vietnam, Argentina, Thailand, Taiwan, and China) have grown faster.

6.8 Long-term growth

Figure 10 shows the number of allocated addresses since 2003 (from RIR data), routed and allocated addresses since 2008 (from RV) and the number of pingable, observed and estimated addresses over time for our 12-month time windows. The ping data from 2003 to 2011 is from USC/LANDER [3]. The ping data since 2012 is the 12-months windows of IPING. Note that allocated and routed addresses are plotted on a different scale (right y-axis).

Two distinct phases are visible for the number of allocated addresses: the last boom 2004–2011, and the slowdown due to running out of unallocated addresses since 2011. The number of allocated addresses increased much faster than the number of pingable addresses until 2011. Even since 2011 the number of allocated

addresses increased faster than the number of pingable addresses. On the other hand, the number of estimated used addresses is increasing much faster than the number of pingable addresses with a growth rate similar to the rate of the allocated and routed addresses before their slowdown in 2011 and 2012 respectively.

6.9 Comparison with Internet user growth

We think the growth of the number of used IPv4 addresses is primarily driven by an Internet user population increase, irrespective of the number of devices per user. All home devices are behind NATs and mobile devices are also largely behind NATs. Similarly, if we look at increasingly complex commercial networks, these are also mainly behind NATs (or not even connected to the Internet). In this section we derive a very rough estimate of the IPv4 address growth based on the Internet user growth and compare it with the growth estimated with CR.

According to data from the ITU [27] the number of Internet users has grown from 16 million in December 1995 to 2.75 billion (roughly 39% of the world’s population) in December 2013 as shown in Figure 11. The growth rate of Internet users looks exponential at the beginning of the graph, however since 2006/2007 the growth appears roughly linear. This is consistent with the roughly linear trend in our results shown in Section 6.3.

Between 2007 and 2012 the number of Internet users grew by roughly 250 million per year (c.f. Figure 11). For private use typ-

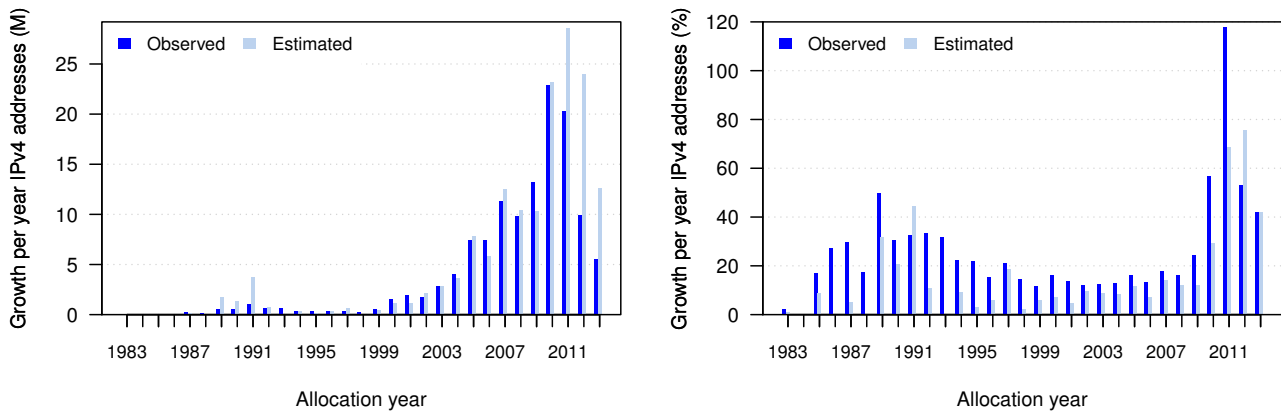


Figure 8: Average absolute and relative yearly growth of observed and estimated IPv4 addresses for different allocation ages

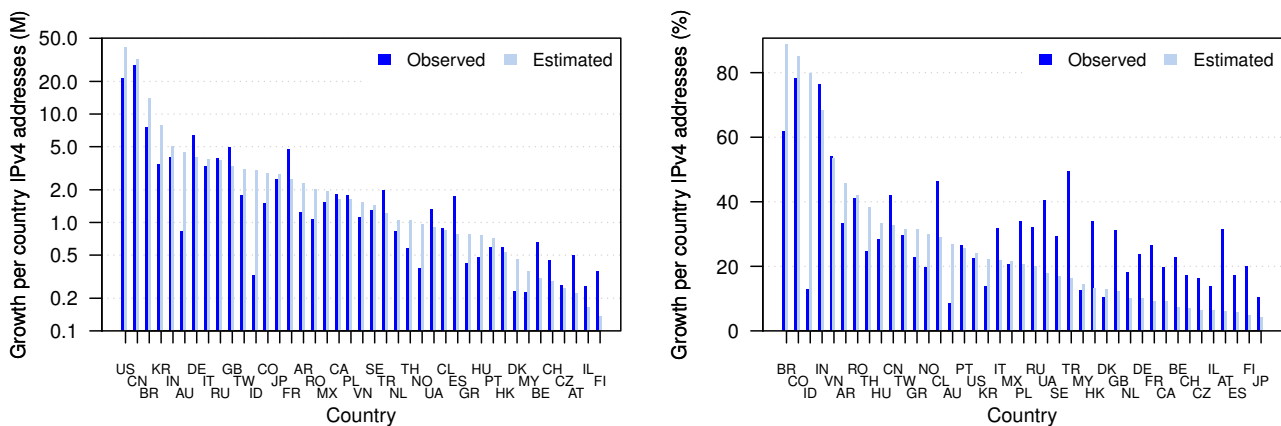


Figure 9: Average absolute and relative yearly growth of observed and estimated IPv4 addresses for different countries sorted by estimated growth (only the largest 42 countries). Note, the absolute numbers are plotted in log scale.

ically a household shares one public IP address. In industrial nations the household size is 2–3, but in developing countries it can be higher, for example it is over 5 in India [28]. We assume the average household size of new Internet users is between 2 and 5.

In addition a fraction of new users will get a public IPv4 address at work. We assume an average employment ratio of 65% [29]. We could not find any data on the number of average public IP addresses used per employee, so we are assuming a wide possible range. As upper limit we assume on average one IPv4 address per two employees, as in reality many employees (especially in developing countries) have no Internet at work, work at home, or share computers with other workers. As a lower limit we assume on average there is one public address per 200 workers.

Let g_U be the user growth per year, H the average household size, p_E the employment ratio, and W the average number of employees sharing an IP at work. Then the IP address growth is $g_I = (1/H + p_E/W) g_U$. Based on the above ranges of H and W , we would expect the IPv4 addresses to grow between 50 million and 205 million per year (plus additional addresses for service and infrastructure growth). Our CR growth estimate from Section 6.3 is 170 million IPv4 addresses per year, which fits within this range.

7. UNUSED SPACE PREDICTION

Our CR technique tells us how many unobserved IPv4 addresses to expect, but says nothing about the distribution of free blocks/prefixes. This is a challenging issue – recipients of newly assigned IPv4 address blocks typically prefer usable-sized contiguous allocations, and forwarding information base (FIB) tables in routers are more efficiently packed if address blocks are allocated hierarchically.

Some information is given by the CR estimate of the used but unobserved /24 networks (in Figure 4). However, this again does not tell us whether these small blocks are isolated or parts of unused larger blocks. In this section, we try to understand how the unseen addresses are distributed among seemingly empty subnets, by observing what happens when data sources are combined sequentially; each new source brings addresses that were unseen by the previous sources, and we can model how those addresses fill the previously empty space.

7.1 Model

Let x_i be the number of observed free / i blocks. Let $Z_{00\dots 0}$ be the total number of new addresses to allocate (given by CR). Let N_i be the number of new addresses assumed to be allocated to vacant / i blocks, assuming sequential allocation. Specifically, if two ad-

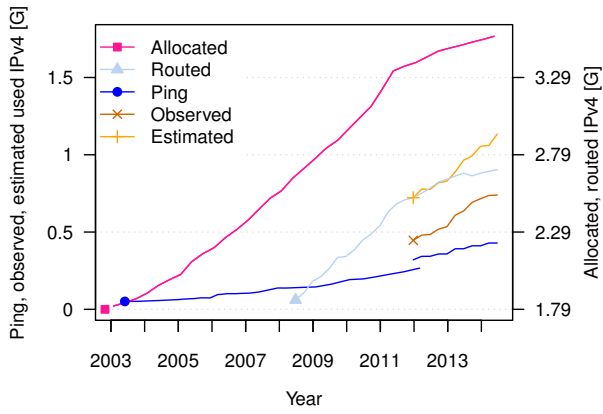


Figure 10: Number of allocated and routed IPv4 addresses (right y-axis), as well as pingable, observed and estimated used IPv4 addresses (left y-axis) over time

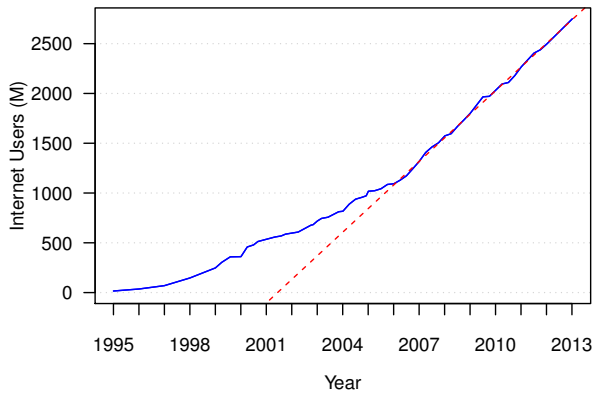


Figure 11: Number of Internet users based on data from ITU

addresses are added to the same vacant $/i$, then only the first of these contributes to N_i , since the block is no longer vacant when the second one is added.

Similarly, let x_i^S be the number of free $/i$ blocks in a set S of addresses, $Z_{00,\dots,0}^{S,\Delta}$ be the number of new addresses when a new set Δ is merged with S , and $n_i^{S,\Delta}$ be the number of addresses added to vacant $/i$ blocks in the process. Without the subscript i , the variables x and n denote vectors.

Note that adding an address to a vacant $/i$ will reduce the number of vacant $/i$ blocks by 1, but increase by one the number of $/j$ blocks for each $j > i$, regardless of where within the $/i$ the address is added. That is,

$$x^{S \cup \Delta} - x^S = An^{S,\Delta}, \quad (2)$$

where

$$A = \begin{pmatrix} -1 & 1 & 1 & \dots & 1 \\ 0 & -1 & 1 & \dots & 1 \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & 0 & \dots & -1 \end{pmatrix}. \quad (3)$$

A natural approach is to estimate the previous fraction of addresses revealed by each new source that have been allocated to free blocks of a given size, and assume the new $Z_{0,0,\dots,0}$ addresses will be distributed in the same way. This is not sufficient, because the allocation process changes the number of available blocks. Instead, our model uses the observation that the probability that a

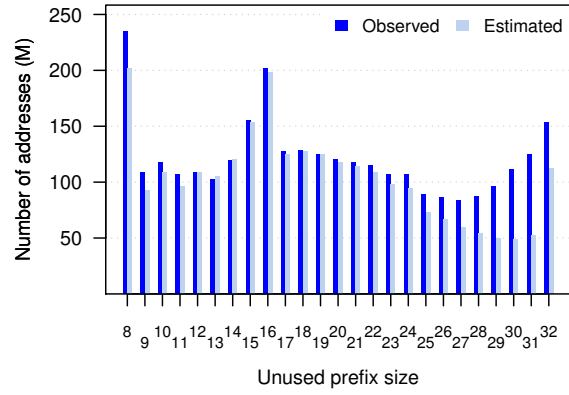


Figure 12: Number of addresses in observed and estimated unused prefixes for different routed prefix sizes

new address is allocated to a free $/i$ block is proportional to x_i , the number of such free blocks. In particular, it assumes that there are f_1, \dots, f_{32} such that the ratio

$$\frac{N_1}{x_1} : \frac{N_2}{x_2 + N_1} : \dots : \frac{N_{32}}{x_{32} + \sum_{j=1}^{31} N_j} = f_1 : f_2 : \dots : f_{32} \quad (4)$$

remains approximately constant as more batches of addresses are discovered.

The model includes subnets larger than $/8$, even though blocks larger than $/8$ have not been allocated. Similarly, we consider all vacant subnets down to vacant $/32$ s, even though subnets smaller than $/24$ are not routed on the public Internet. However, before computing the remaining unused prefixes we split a few $/7$ into $/8$, and we also exclude all private, multicast, experimental and reserved prefixes, such as $224.0.0.0/3$ or $10.0.0.0/8$. Note that we do not exclude non-publicly routed prefixes.

To determine how the unobserved addresses predicted by CR will affect the distribution of free blocks, it remains to determine f_i . To do this, we observe the change in x when a new data set Δ is added to an existing list S of used IPv4 addresses, and from that calculate n . Since A in (3) is invertible, (2) gives

$$n^{S,\Delta} = A^{-1}(x^{S \cup \Delta} - x^S).$$

The f_i are then found by (4), normalized so that $f_{32} = 1$. Since few large subnets become newly used for each data set, estimates of f_i for $i \lesssim 12$ are noisy. This is unfortunate, since these are the blocks of greatest interest. To reduce this noise, estimates were averaged over four cases: $\Delta = \text{IPING, GAME, WEB, WIKI}$; in each case, S is the union of all remaining datasets, except SWIN and CALT.

One concern with this model is that, as the address space fills up, the values of f_i may vary. To check this, we performed tests where datasets were added to S one at a time, in both increasing and decreasing order of the dataset size. The values were reasonably consistent in each case.

7.2 Results

Figure 12 shows the number of addresses in unused prefixes at 30 June 2014, based on all sources except SWIN and CALT. Results are for both direct observation and CR. The majority of empty prefixes are longer than $/20$ (fewer than 2^{12} addresses), but the unused space is roughly uniformly distributed among prefixes of lengths $/9$ to $/24$ (except $/15$ and $/16$). The reason for this is unclear.

If the used but unobserved $/8$ to $/24$ subnets estimated by the model of Section 7.1 were divided into $/24$ s, there would be 0.3

million /24s. This is consistent with the estimate of 0.26–0.36 million by the independent LLM model for the period ending at 30 June 2014, providing evidence for the validity of both models.

7.2.1 Router FIB limitations

One of the reasons that the distribution of prefix sizes is important is that each routed prefix requires an entry in a router’s Forwarding Information Base (FIB), and there is a *prima facie* risk that allocating all unused prefixes could overflow the FIBs. Above we estimated that including the unrouted space there are 0.78 million prefixes that are /24 or larger. Currently, there are more than 0.5 million routed prefixes already (but a substantial fraction is unused). In 2007 Juniper [30] stated that its M120 and MX960 had FIB capacities of about 2 million IPv4 routes, and that IPv4 FIBs with approximately 10 million entries are feasible within a few years if demand exists. In addition, FIB compression techniques can reduce size of FIBs [30]. This suggests that it will be feasible to use and route all less than 1.3 million available prefixes. Even if unused prefixes are subdivided further, it appears feasible to route them all. However, some existing routers may require upgrades.

7.2.2 Estimated years of supply

In July 2014 roughly 5.5 /8 networks of unallocated addresses remained [1] – equivalent to 350,000 /24 networks or 90 million addresses. Even if in addition to these all 4.4 million routed but unused /24 subnets could be prised away from their current owners, /24 networks would be exhausted in 2024 under the current growth trend of 0.45 million /24 subnets per year. Unused routed IPv4 addresses would be exhausted in 2023 given the current growth of 170 million addresses per year.

The overall estimate hides the differences between regions. Table 6 shows the available space (unallocated space plus allocated publicly routed unused space based on our CR estimates), the current average growth rate over our measured time period and the predicted year supply will run out (under the very optimistic assumption that the *whole* unused space could be utilised) for IPs and /24 networks for each RIR. For most RIRs the number of years of IP supply is equal or smaller than the number of years of /24 networks supply. This is because either the IP and /24 growth rates are similar or the IP growth rate is larger, but there are significant numbers of unused IPs in used /24. RIPE is an exception due to similar growth rates for IPs and /24s and significant supply in already used blocks. At current growth rates ARIN and RIPE have 14+ years of IP supply left, but AfriNIC has only 8–9 years of supply left, and LACNIC and APNIC have only 2–4 years of supply left. Any future reallocations between RIRs to ease local pressures would change these numbers of course.

However, it appears unlikely that the whole IPv4 space will ever be completely utilised. If the overall utilization of routed /24 subnets remains below, say 75%, the current growth rates suggests four years of remaining supply overall. Moreover, in this case regions with tight supply, such as APNIC and LACNIC, would be exhausted within one year. APNIC stands out, because its IP growth rate is much larger than its /24 growth rate with unused IPs in used /24s depleted soon. We expect a slowdown shortly – in fact Figure 6 suggests it may have begun already. An open question is the large amount of unrouted IPv4 space, much of which has not been routed for years. Unused parts of the unrouted space might provide a short-lived increase in IPv4 supply.

Over the next one to two years, we expect IPv4 exhaustion to be increasingly felt, resulting in a brief growth in the IPv4 address market. Most organizations holding unused addresses do so for operational reasons – to allow expansion or flexibility, or in one

case as a /8 darknet – but some may be holding them to sell if the market price rises sufficiently. The numbers in this paper may guide how long they can expect to be held for, assuming that the market will collapse once IPv6 is widely adopted. However, this is complicated by the fact that the very act of selling a large block of IPv4 addresses will delay the implementation of IPv6, and hence prolong the IPv4 market.

8. CONCLUSIONS AND FUTURE WORK

Our key contribution is describing and demonstrating a new statistical *capture-recapture* technique for improved estimation of the true population of both observed and unobserved (yet still active) IPv4 addresses from diverse sources of active and passive measurement data. This technique refines our community’s understanding of IPv4 address space exhaustion and consequent incentives for IPv6 adoption.

Data from nine sources over the past three years suggests 5.9 million used /24 subnets and 740 million used IPv4 addresses. Yet our CR technique indicates a significantly higher 1.2 billion IPv4 addresses in use across 6.3 million /24 subnets, with usage growing at around 0.45 million /24 subnets and 170 million IPv4 addresses per year. Asia and Europe have the highest numbers of used IP addresses, while Africa and South America show the fastest growth. Based on the overall estimates, at the very best unallocated plus routed but unused addresses will last until 2023. More likely, if say only 75% of routed /24 subnets could ever be used, supply will be exhausted in 2018. Moreover, Asia and South America are the two “pressure points” that will experience a shortage of addresses within the next one or two years.

Previous sales of IP addresses saw prices of US\$8–17 per IP depending on the sizes of the blocks sold [31, 32]. At an average price of US\$10 per IP address, the 4.4 million routed unused /24 subnets have a value of over US\$11 billion. However, probably only a small fraction of those will be sold, so even if prices rise substantially, the eventual market value is likely to be smaller.

Our ability to collect more IP data for validating or improving our estimates, or potentially detecting more hosts (e.g. private servers), is limited by common privacy restrictions. We plan to explore an enhanced method [33] for securely applying CR to multi-source measurement data without revealing which IPv4 addresses each source contains.

Acknowledgements

This research was supported by Australian Research Council grants LP110100240 (with APNIC Pty Ltd) and FT0991594. We thank Geoff Huston, George Michaelson, Valve Corporation, A. Reynolds, Swinburne ITS, Caltech IMSS, D. Buttigieg, C. Tassios, R. Bevier, B. Mattern, USC/ISI and J. Heidemann for providing data. We thank our shepherd X. Dimitropoulos and the anonymous reviewers for their helpful comments.

9. REFERENCES

- [1] G. Huston. IPv4 Address Report. <http://www.potaroo.net/tools/ipv4/index.html>.
- [2] Y. Pryadkin, R. Lindell, J. Bannister, R. Govindan. An Empirical Evaluation of IP Address Space Occupancy. Technical Report ISI-TR 598, USC/ISI, 2004.
- [3] J. Heidemann, Y. Pradkin, R. Govindan, C. Papadopoulos, G. Bartlett, J. Bannister. Census and Survey of the Visible Internet. In *ACM Conference on Internet measurement (IMC)*, pages 169–182, 2008.

Table 6: Available IPv4 addresses and /24 networks (unallocated plus publicly routed but unused), growth rate and years where supply will run out (under the very optimistic assumption that all currently unused space will be used) by RIR

RIR	Available IPs [M]	Growth IPs [M/year]	Year runout IPs	Available /24s [M]	Growth /24s [M/year]	Year runout /24s
AfriNIC	75	9	2022–2023	0.29	0.03	2023–2024
APNIC	310	80	2017–2018	1.04	0.11	2023–2024
ARIN	830	35	2037–2038	2.39	0.08	2043–2044
LACNIC	65	25	2016–2017	0.15	0.09	2015–2016
RIPE	370	25	2028–2029	0.87	0.11	2021–2022
World	1650	170	2023–2024	4.7	0.45	2024–2025

- [4] X. Cai, J. Heidemann. Understanding Block-level Address Usage in the Visible Internet. In *ACM SIGCOMM Conference*, pages 99–110, 2010.
- [5] Internet Census 2012 – Port scanning /0 using insecure embedded devices, 2012. <http://internetcensus2012.bitbucket.org>.
- [6] A. Dainotti, K. Benson, A. King, kc claffy, M. Kallitsis, E. Glatz, X. Dimitropoulos. Estimating Internet Address Space Usage Through Passive Measurements. *ACM Computer Communication Review (CCR)*, 44(1):42–49, Jan. 2014.
- [7] C. G. J. Petersen. The Yearly Immigration of Young Plaice into the Limfjord from the German Sea. *Rept. Danish Biol. Sta.*, 6:1–77, 1895.
- [8] F. C. Lincoln. Calculating Waterfowl Abundance on the Basis of Banding Returns. *U.S. Dept. Agric. Circ.*, 118:1–4, 1930.
- [9] A. Chao. An Overview of Closed Capture-Recapture Models. *Journal of Agricultural, Biological, and Environmental Statistics*, 6(2):158–175, 2001.
- [10] S. Zander, L. L. H. Andrew, G. Armitage, and G. Huston. Estimating IPv4 Address Space Usage with Capture-Recapture. In *7th IEEE Workshop on Network Measurements (WNM)*, Oct. 2013.
- [11] S. Zander, L. L. H. Andrew, G. Armitage, G. Huston, and G. Michaelson. Mitigating Sampling Error when Measuring Internet Client IPv6 Capabilities. In *ACM Internet Measurement Conference (IMC)*, Nov. 2012.
- [12] DNS-based Blacklist of NiX Spam. <http://www.dnsbl.manitu.net/>.
- [13] Measurement Lab. <http://www.measurementlab.net/>.
- [14] University of Oregon Route Views Project. <http://www.routeviews.org/>.
- [15] X. Meng, Z. Xu, B. Zhang, G. Huston, S. Lu, L. Zhang. IPv4 Address Allocation and the BGP Routing Table Evolution. *ACM Computer Communication Review (CCR)*, 35(1):71–80, 2005.
- [16] A. Sriraman, K. R. B. Butler, P. D. McDaniel, P. Raghavan. Analysis of the IPv4 Address Space Delegation Structure. In *IEEE Symposium on Computers and Communications (ISCC)*, pages 501–508, Jul. 2007.
- [17] E. B. Hook, R. R. Regal. Capture-Recapture Methods in Epidemiology: Methods and Limitations. *Epidemiologic Reviews*, 17(2):243–264, 1995.
- [18] M. Roughan, J. Tuke, O. Maennel. Bigfoot, Sasquatch, the Yeti and other missing links: what we don’t know about the AS graph. In *8th ACM Internet Measurement Conference (IMC)*, pages 325–330, Oct. 2008.
- [19] A. Chao, P. K. Tsay, S. H. Lin, W. Y. Shau, D. Y. Chao. The Applications of Capture-Recapture Models to Epidemiological Data. *Statistics in Medicine*, 20:3123–3157, Oct. 2001.
- [20] S. E. Fienberg. The Multiple Recapture Census for Closed Populations and Incomplete 2k Contingency Tables. *Biometrika*, 59(3):591–603, Dec. 1972.
- [21] E. Cooch, G. C. White. *Program MARK: A Gentle Introduction*. Cornell University, 2009.
- [22] K. P. Burnham, D. R. Anderson. Multimodel Inference - Understanding AIC and BIC in Model Selection. *Sociological Methods & Research*, 33:261–304, 2004.
- [23] S. Baillargeon, L.-P. Rivest. Rcapture: Loglinear Models for Capture-Recapture in R. *Journal of Statistical Software*, 19(5):1–31, Apr. 2007.
- [24] Internet Addresses Census dataset, PREDICT ID: USC-LANDER/internet_address_census_it40c-20110406. Provided by the USC/LANDER project. <http://www.isi.edu/ant/lander>.
- [25] D. Moore, G. M. Voelker, S. Savage. Inferring Internet Denial-of-Service Activity. In *Usenix Security Symposium*, August 2001.
- [26] E. Gregori, A. Improta, L. Lenzi, L. Rossi, L. Sani. On the Incompleteness of the AS-level Graph: a Novel Methodology for BGP Route Collector Placement. In *Internet Measurement Conference (IMC)*, 2012.
- [27] ITU key 2006-2013 ICT data for the world, 2013. http://www.itu.int/en/ITU-D/Statistics/Documents/statistics/2013/ITU_Key_2005-2013_ICT_data.xls.
- [28] Wikipedia. List of countries by number of households, Oct. 2013. http://en.wikipedia.org/w/index.php?title=List_of_countries_by_number_of_households&oldid=576467223.
- [29] Wikipedia. Employment-to-population ratio, Mar. 2014. http://en.wikipedia.org/w/index.php?title=Employment-to-population_ratio&oldid=598945003.
- [30] J. Scudder. Router Scaling Trends. Presentation at RIPE-54, May 2007. http://meetings.ripe.net/ripe-54/presentations/Router_Scaling_Trends.pdf.
- [31] B. Edelman and M. Schwarz. Pricing and Efficiency in the Market for IP Addresses. Working Paper Number: 12-020, Nov. 2011. <http://hbswk.hbs.edu/item/6849.html>.
- [32] S. Brown. IPv4 Trading in Review, Jan. 2014. <http://ipv4marketgroup.com/blog/>.
- [33] S. Zander, L. L. H. Andrew, and G. Armitage. Estimating the used IPv4 address space with secure multi-party capture-recapture. In *INFOCOM (poster)*, Turin, Italy, 15-18 Apr 2013.