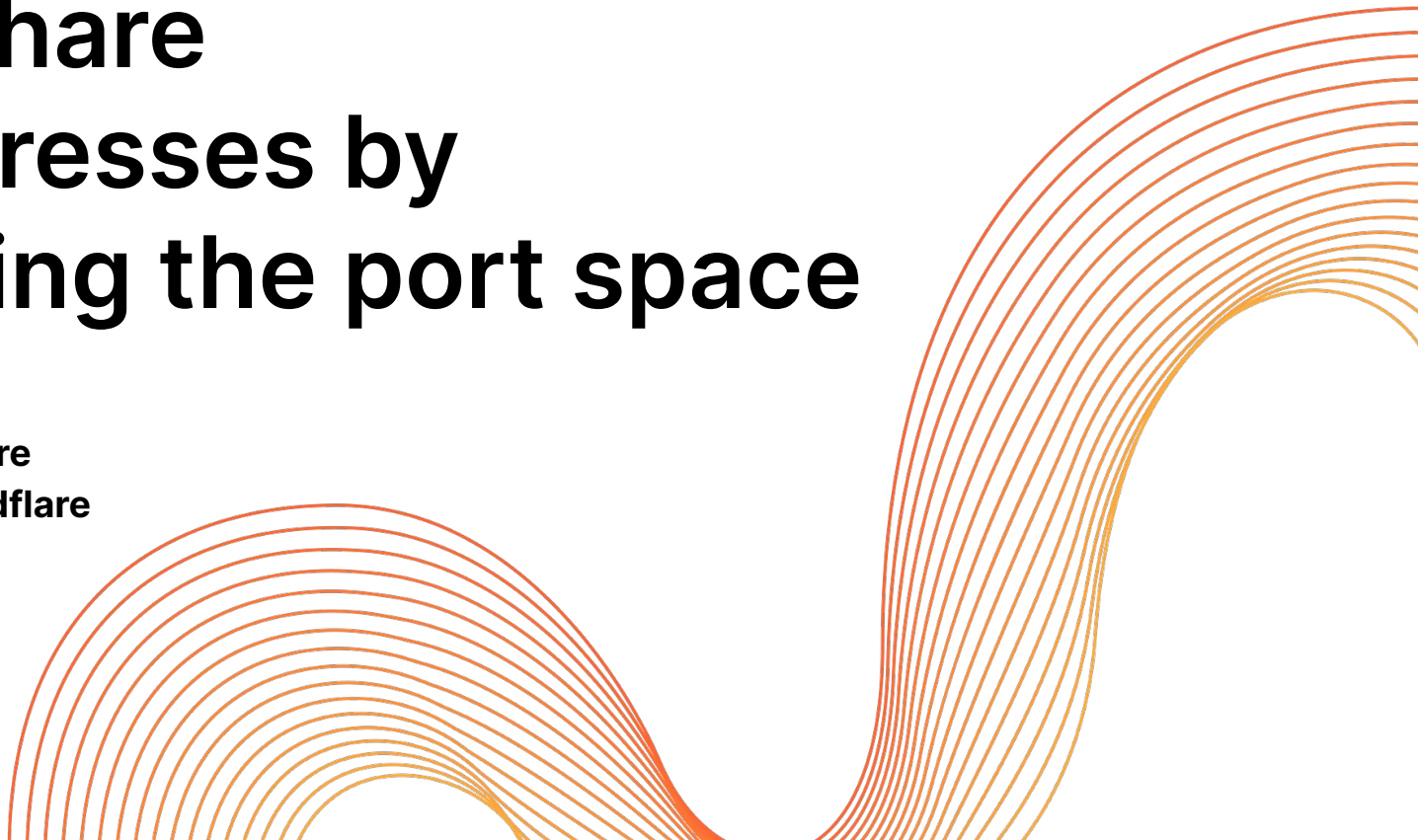


How to share IPv4 addresses by partitioning the port space

Jakub Sitnicki, Cloudflare
Marek Majkowski, Cloudflare

Linux Plumbers Conference
Dublin, Ireland
September 12-14, 2022



Part I

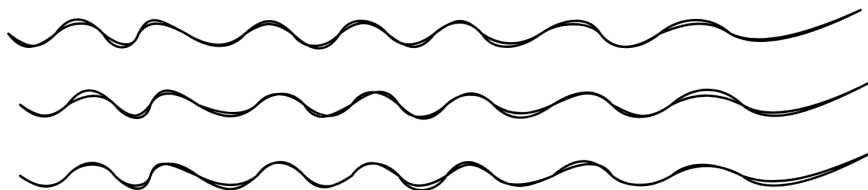
The Situation

Not all public IP addresses are created equal



198.51.100.183 != 203.0.113.72

198.51.100.183 -> geo-location

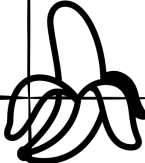
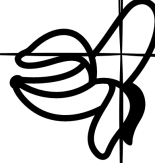
 

1 2 3 4 ...

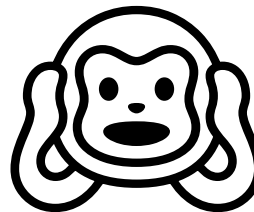
-
- Dublin, Ireland - From your IP address

198.51.100.183 -> reputation

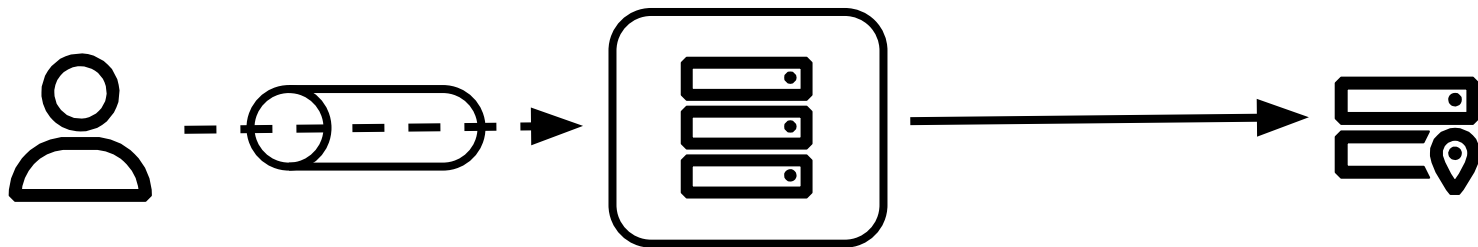
Select all squares with
bananas

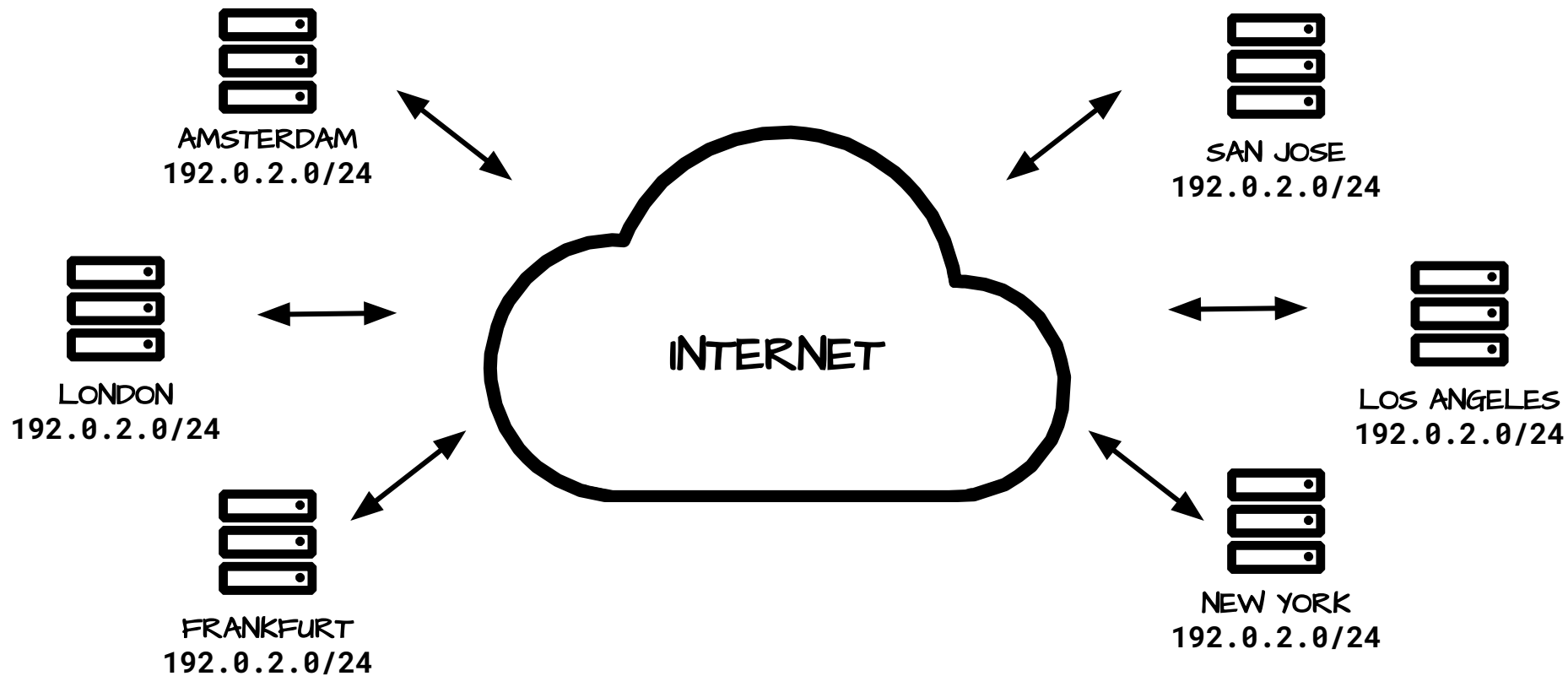
VERIFY



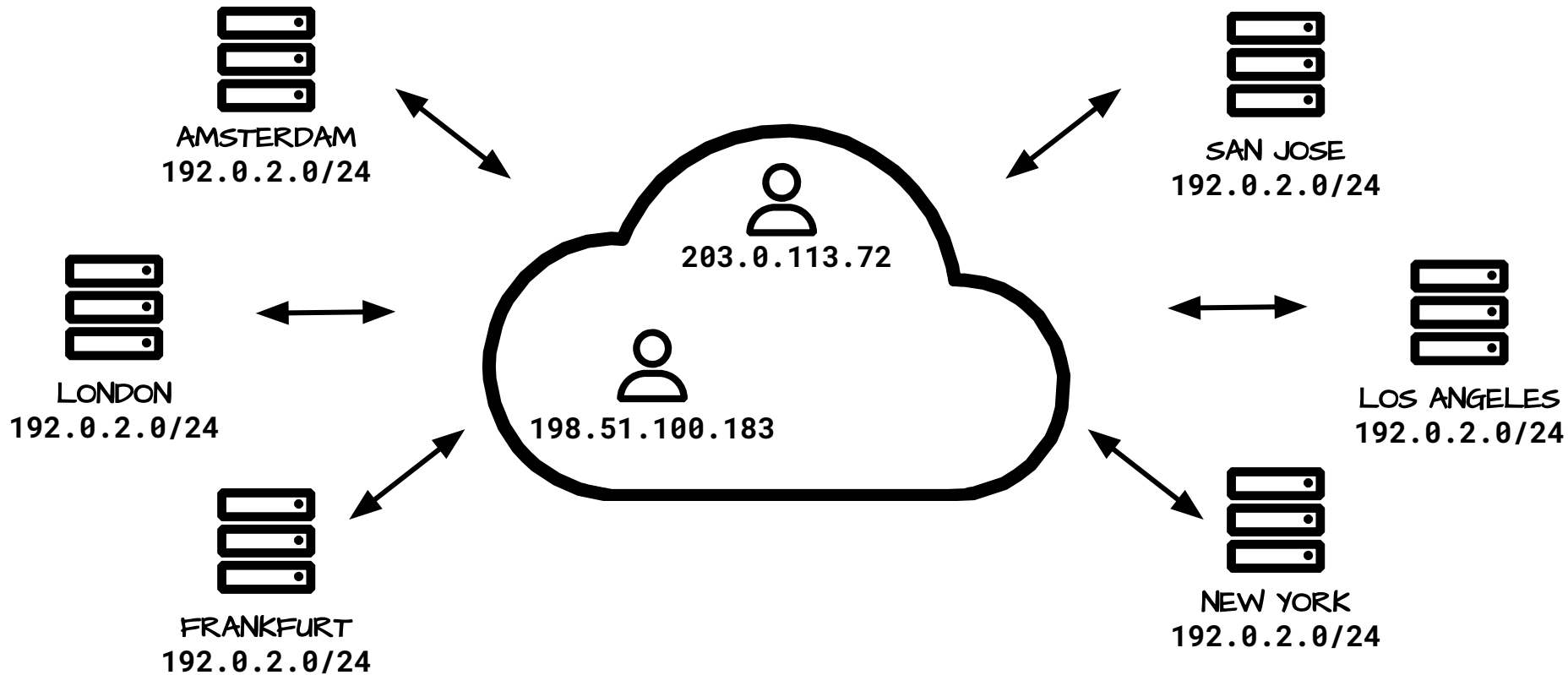
We run a VPN called WARP ©



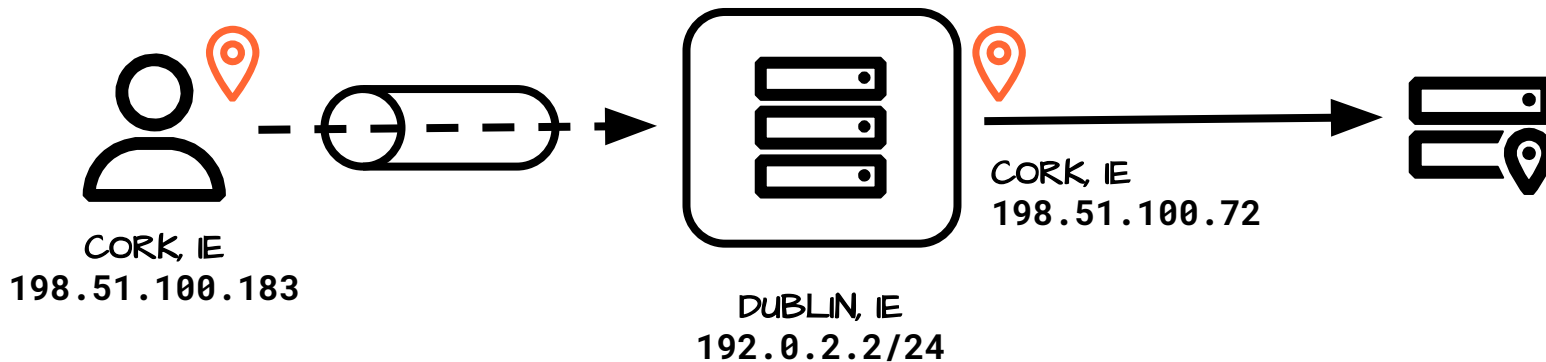
We run an anycast network



We run an anycast network



 -> (location, reputation, ...)



Part II

The Problem

Cork IP on every CF server

Need lots of IPs

Egress always local

Excellent reliability

Cork IP on one server

Need just one IP

Forward for egress

Poor reliability

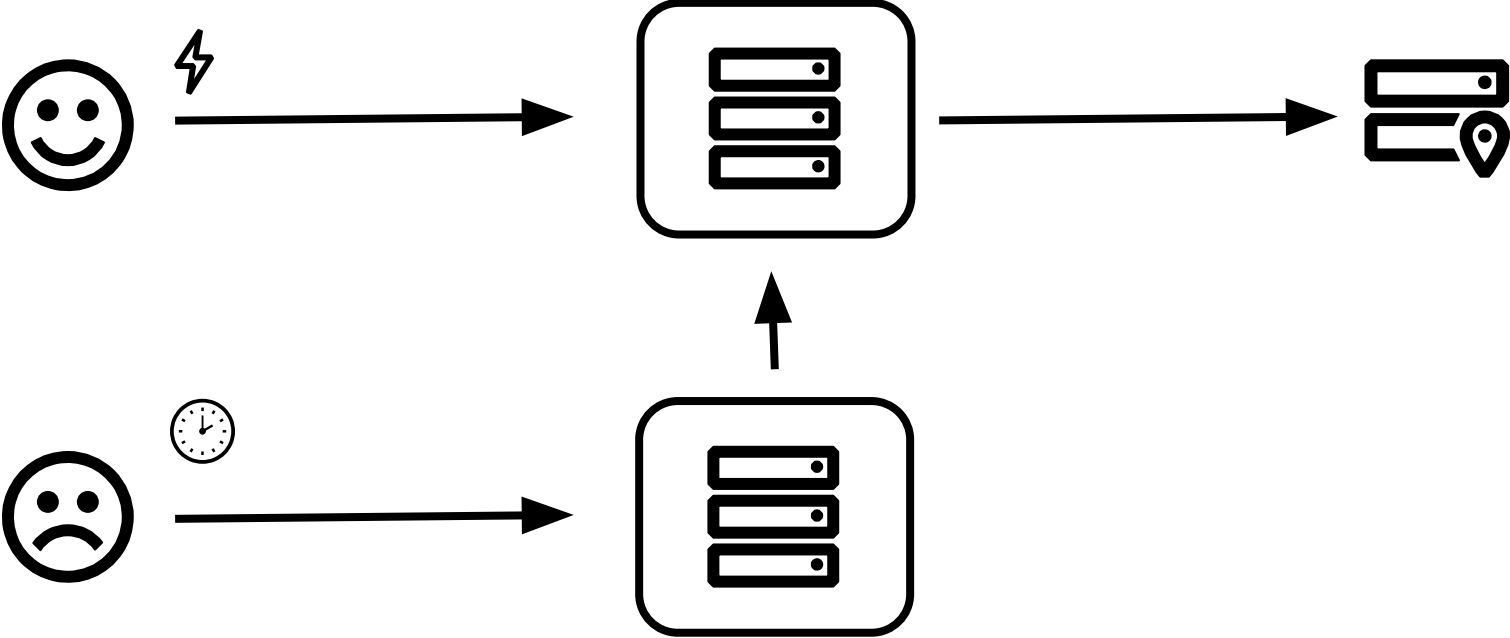
\$\$\$



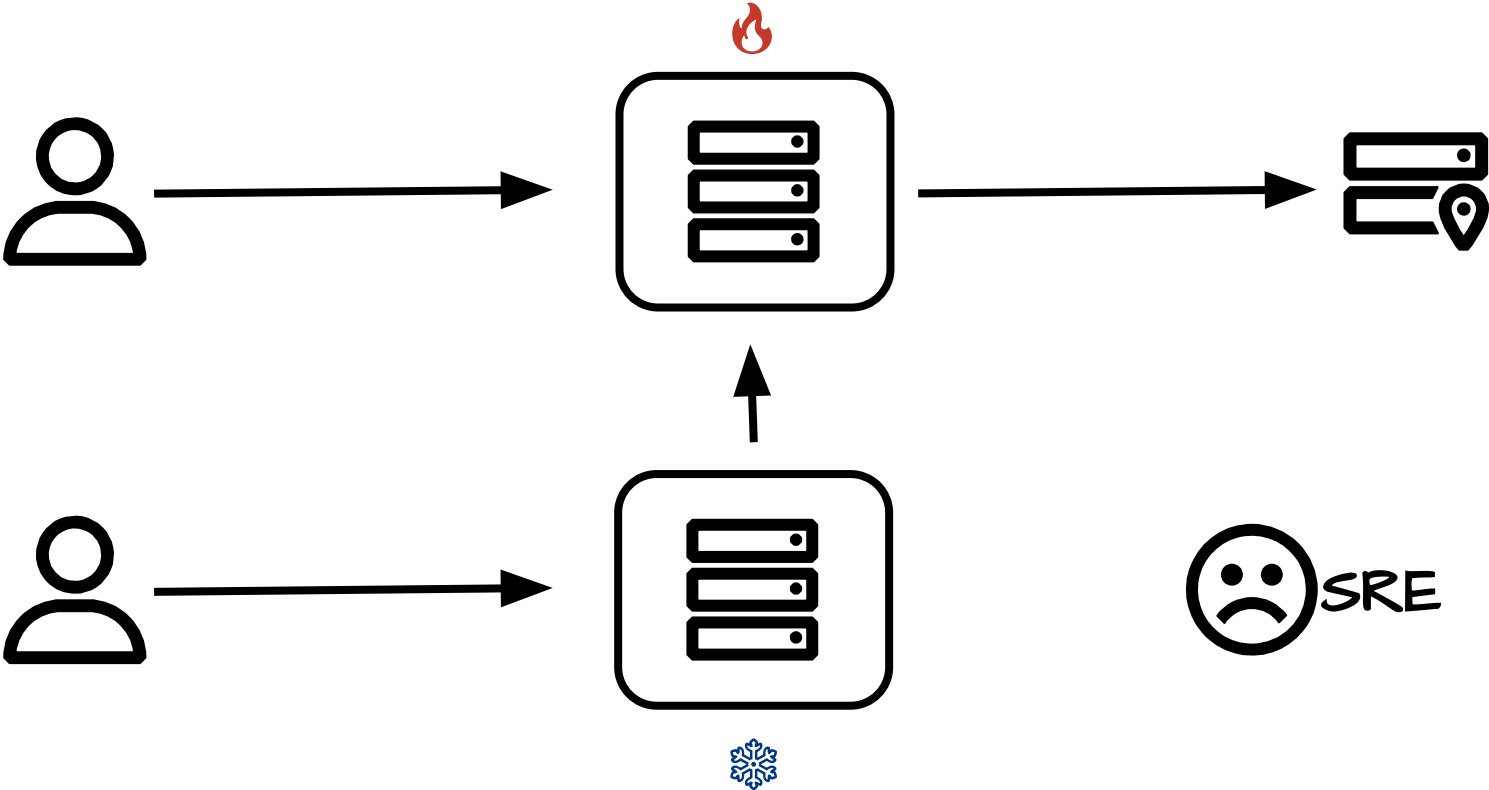
\$

Cork IP on every server	Cork IP per DC	Cork IP on one server
Need lots of IPs	Need some IPs	Need just one IP
Egress always local	Egress mostly local	Forward for egress
Excellent reliability	Good reliability	Poor reliability
\$\$\$	\$\$	\$

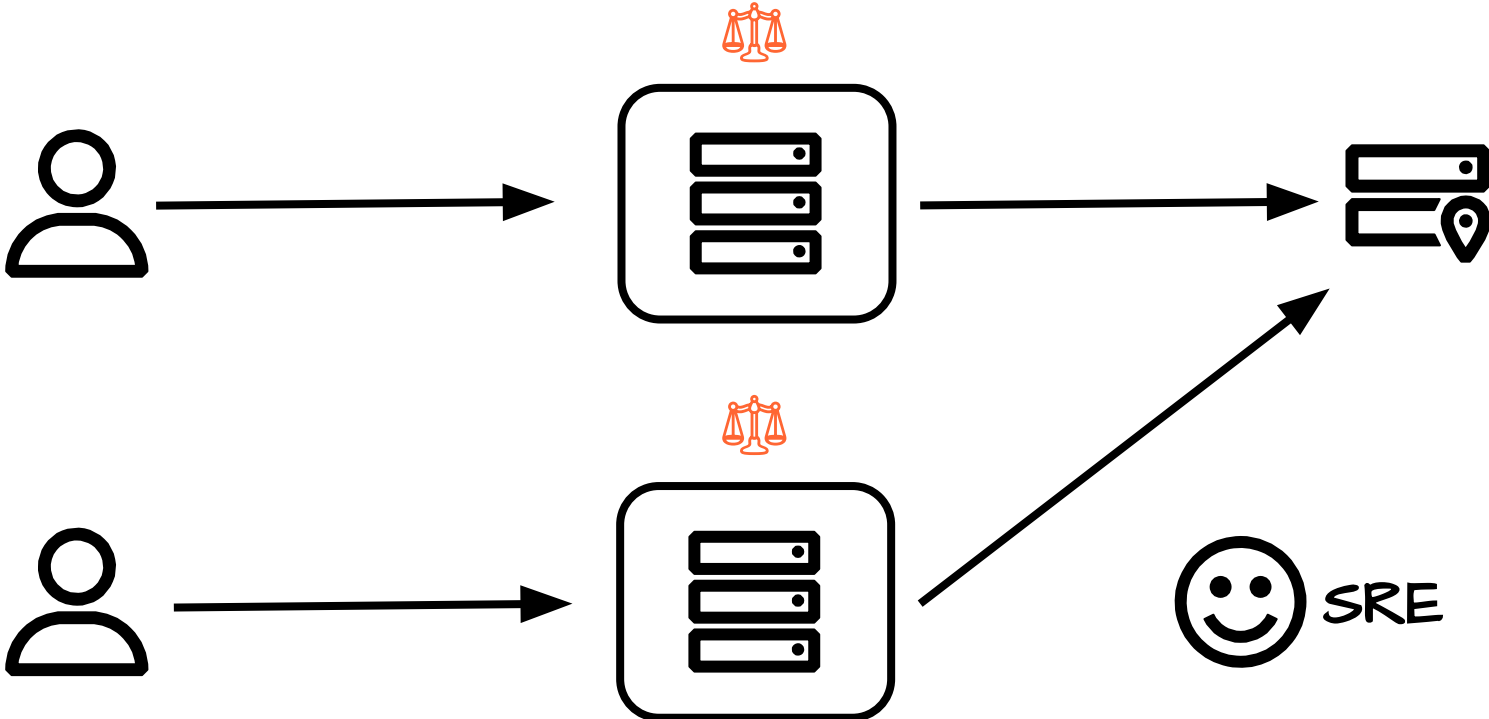
Local egress is happy user



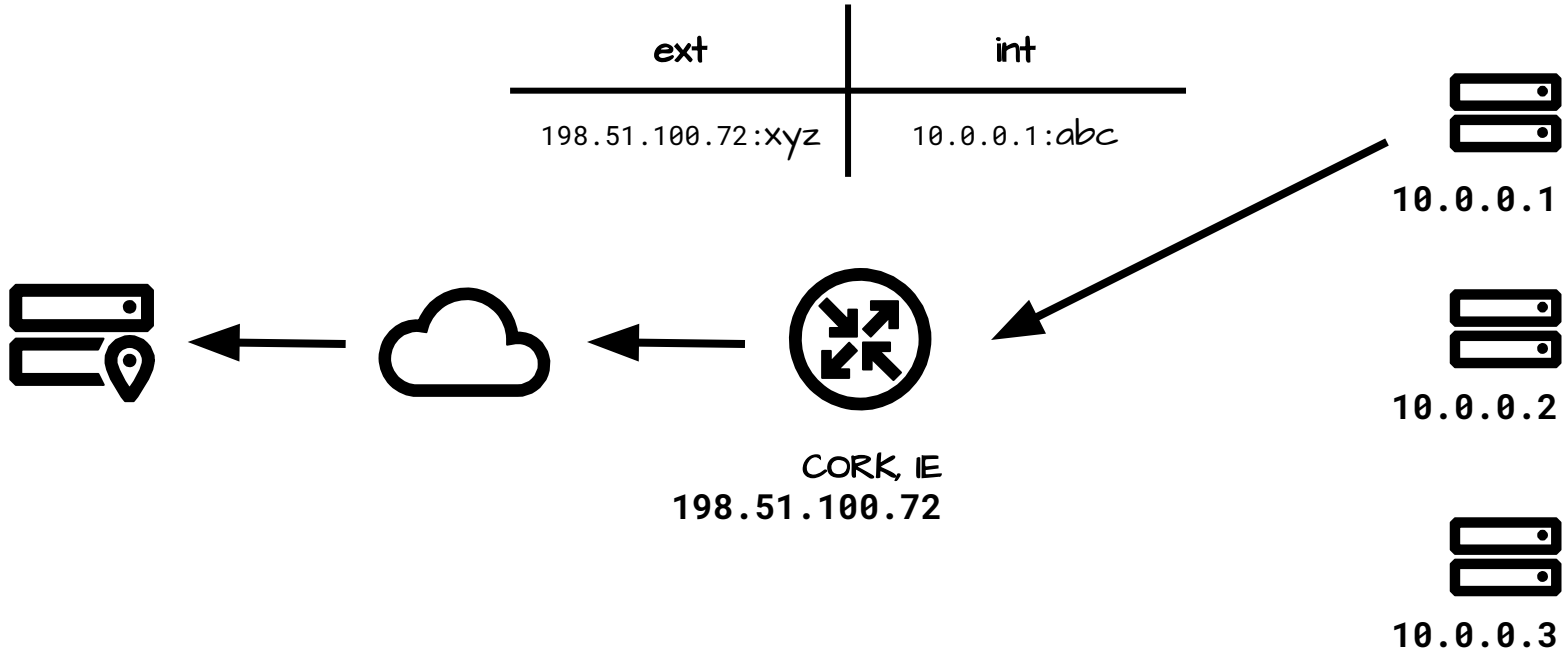
Local egress is happy SRE



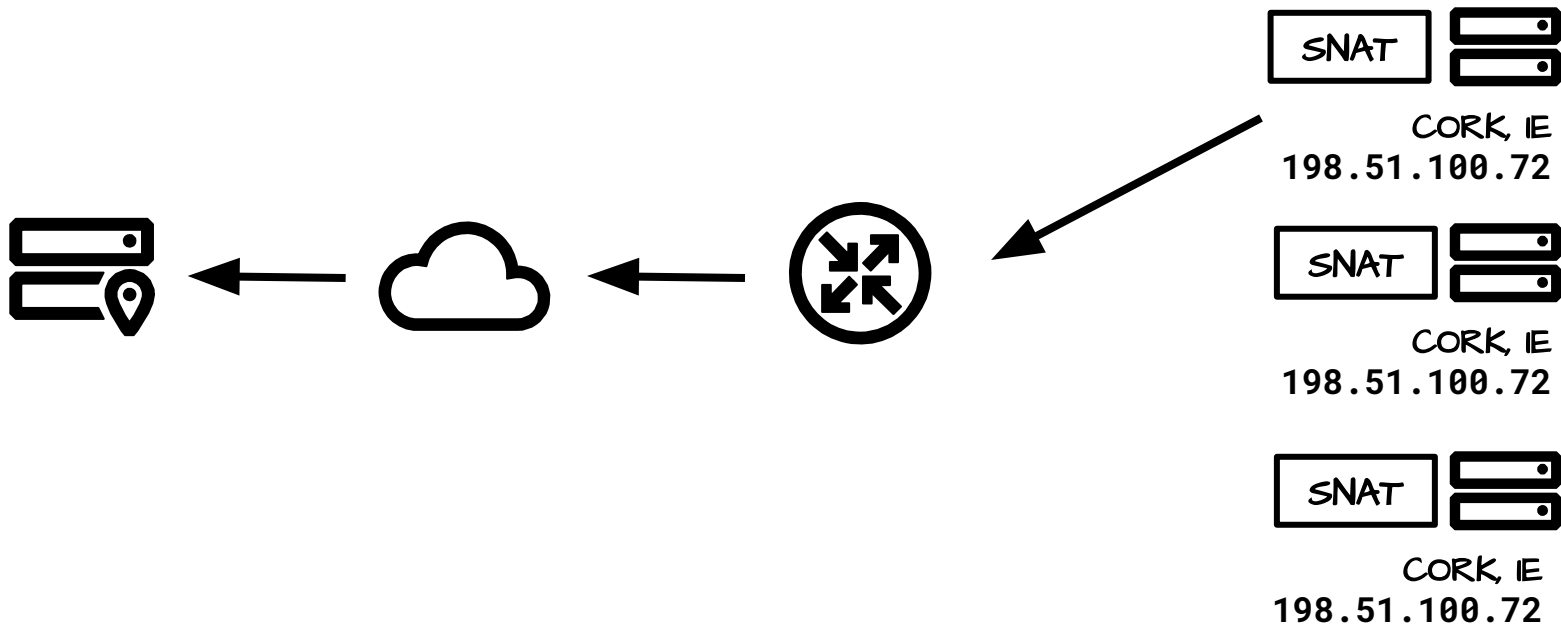
Local egress is happy SRE



SNAT on router?



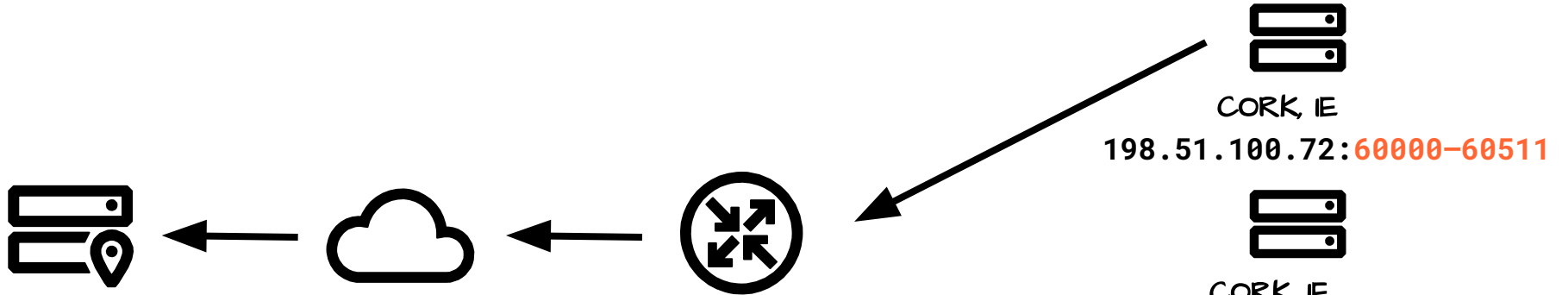
Distributed SNAT?



Part III

The Solution

Share the IP, partition the port range



Each server gets 512 egress ports per IP
Return traffic routed by (dst IP, dst port) !

CORK, IE
198.51.100.72:61024-61535

How to open connections from a given port range?

```
system("sysctl -w net.ipv4.ip_local_port_range='60000 60511'")
```



```
s = socket(AF_INET, SOCK_STREAM)
s.setsockopt(SOL_IP, IP_BIND_ADDRESS_NO_PORT, 1) # share the local port
s.bind(("192.0.2.1", 0)) # let the kernel find a free 4-tuple
s.connect(("1.1.1.1", 53))
```

Sample the port range, try each port

```
{ 60,000 ; 60,002 ; 60,003 ; 60,007 ; 60,009 }
```


```
s1 = socket(AF_INET, SOCK_STREAM)
s1.setsockopt(SOL_SOCKET, SO_REUSEADDR, 1)
s1.bind(("192.0.2.1", 60_000))
s1.connect(("1.1.1.1", 53)) # EADDRNOTAVAIL
```

```
s2 = socket(AF_INET, SOCK_STREAM)
s2.setsockopt(SOL_SOCKET, SO_REUSEADDR, 1)
s2.bind(("192.0.2.1", 60_002))
s2.connect(("1.1.1.1", 53)) # OK
```



4-tuple is busy
try next port



 must use REUSEADDR to share the local (IP, port)

 must re-create the socket each time due to PORTLOCK

Sample the port range, try each port

```
{ 60,000 ; 60,002 ; 60,003 ; 60,007 ; 60,009 }
```

```
s1 = socket(AF_INET, SOCK_STREAM)
s1.setsockopt(SOL_SOCKET, SO_REUSEADDR, 1)
s1.bind(("192.0.2.1", 60_000))
s1.connect(("1.1.1.1", 53))
```

```
s2 = socket(AF_INET, SOCK_STREAM)
s2.setsockopt(SOL_SOCKET, SO_REUSEADDR, 1)
s2.bind(("192.0.2.1", 60_002))
s2.connect(("1.1.1.1", 53))
```

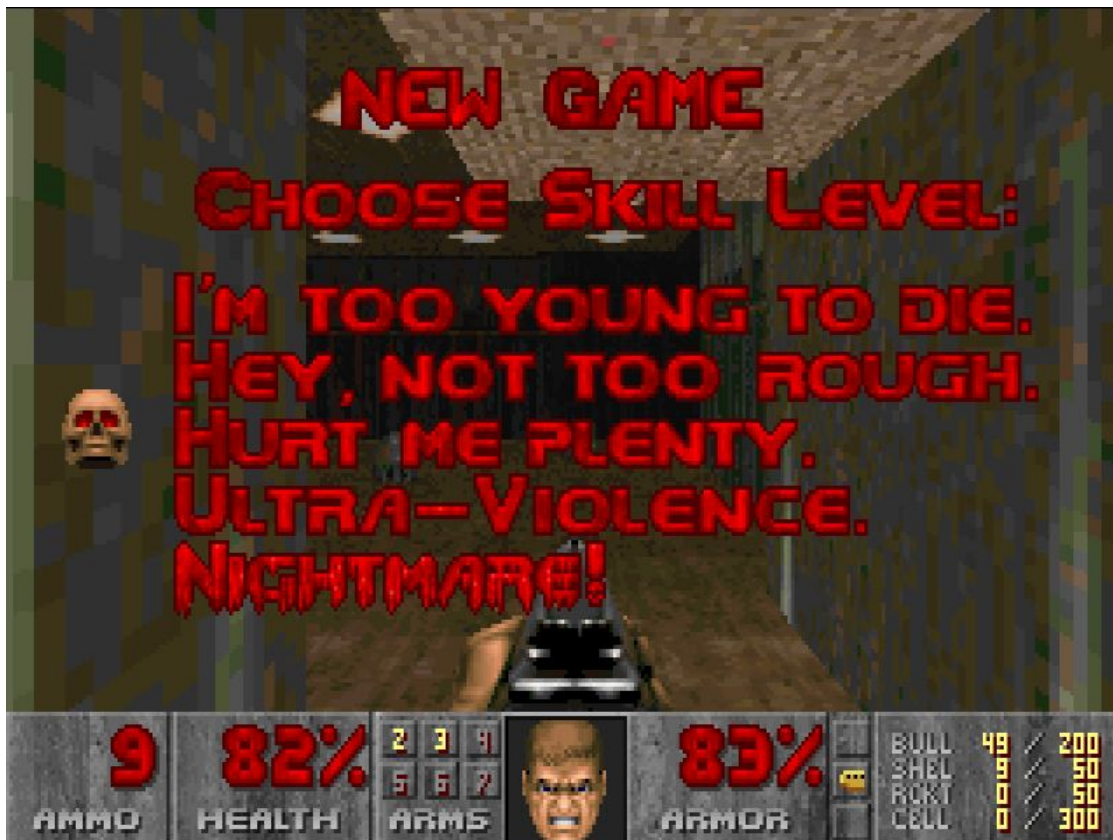
TCP ONLY

4-tuple is busy
try next port

Sample the port range, try each port

What about UDP?

TCP & UDP
semantics differ



TCP & UDP connect () semantics differ

```
s1.setsockopt(SOL_SOCKET, SO_REUSEADDR, 1)
s2.setsockopt(SOL_SOCKET, SO_REUSEADDR, 1)
```

	TCP	UDP
s1.bind(192.0.2.1, 60000)	⇒ OK	⇒ OK
s1.connect(1.1.1.1, 53)	⇒ OK	⇒ OK
s2.bind(192.0.2.1, 60000)	⇒ OK	⇒ OK
s2.connect(1.1.1.1, 53)	⇒ EADDRINUSE	⇒ OK

Must use REUSEADDR to share the local (IP, port)
Otherwise -> only as many concurrent flows as local ports

UDP port conflict detection in user-space

```
s = socket(AF_INET, SOCK_DGRAM)
s.setsockopt(SOL_SOCKET, SO_REUSEADDR, 1)
s.bind(("192.0.2.1", 60_000))

cookie = s.getsockopt(SOL_SOCKET, SO_COOKIE, 8)
s.setsockopt(SOL_SOCKET, SO_REUSEADDR, 0) # block port
c = netlink_udp_lookup(("192.0.2.1", 60_000), ("1.1.1.1", 53))
if c != cookie: # conflict?
    raise OSError(EADDRNOTAVAIL)
s.connect(("1.1.1.1", 53))
s.setsockopt(SOL_SOCKET, SO_REUSEADDR, 1) # unblock port
```

Part IV

The Proposal

Again - How to open connections from a given port range?

```
system("sysctl -w net.ipv4.ip_local_port_range='60000 60511'")
```



```
s = socket(AF_INET, SOCK_STREAM)
s.setsockopt(SOL_IP, IP_BIND_ADDRESS_NO_PORT, 1) # share the local port
s.bind(("192.0.2.1", 0)) # let the kernel find a free 4-tuple
s.connect(("1.1.1.1", 53))
```

Ephemeral port range per socket

```
system("sysctl -w net.ipv4.ip_local_port_range='60000 60009'")
```

```
s = socket(AF_INET, SOCK_STREAM)
s.setsockopt(SOL_IP, IP_BIND_ADDRESS_NO_PORT, 1)
s.setsockopt(SOL_IP, IP_LOCAL_PORT_RANGE, [60000, 60009])
s.bind(("192.0.2.1", 0)) # let the kernel find a free 4-tuple
s.connect(("1.1.1.1", 53))
```

What if we could do this? 

Ephemeral port range per socket

✉ RFC posted to netdev

```
--- a/include/net/inet_sock.h
+++ b/include/net/inet_sock.h
@@ -238,6 +238,10 @@ struct inet_sock {
     __be32                mc_addr;
     struct ip_mc_socklist __rcu    *mc_list;
     struct inet_cork_full  cork;
+
+     struct {
+         __u16 lo;
+         __u16 hi;
+     } local_port_range;
};
```

Set in `ip_sockglue.c` by `IP_LOCAL_PORT_RANGE` option handler

Ephemeral port range per socket

```
void inet_sk_get_local_port_range(const struct sock *sk, int *low, int *high)
{
    const struct inet_sock *inet = inet_sk(sk);
    const struct net *net = sock_net(sk);
    int lo, hi;

    inet_get_local_port_range(net, &lo, &hi);

    if (unlikely(inet->local_port_range.lo))
        lo = clamp_val(inet->local_port_range.lo, lo, hi);
    if (unlikely(inet->local_port_range.hi))
        hi = clamp_val(inet->local_port_range.hi, lo, hi);

    *low = lo;
    *high = hi;
}
```

Called in `->connect()` \Rightarrow `tcp_v4_connect()` \Rightarrow `inet_hash_connect()`

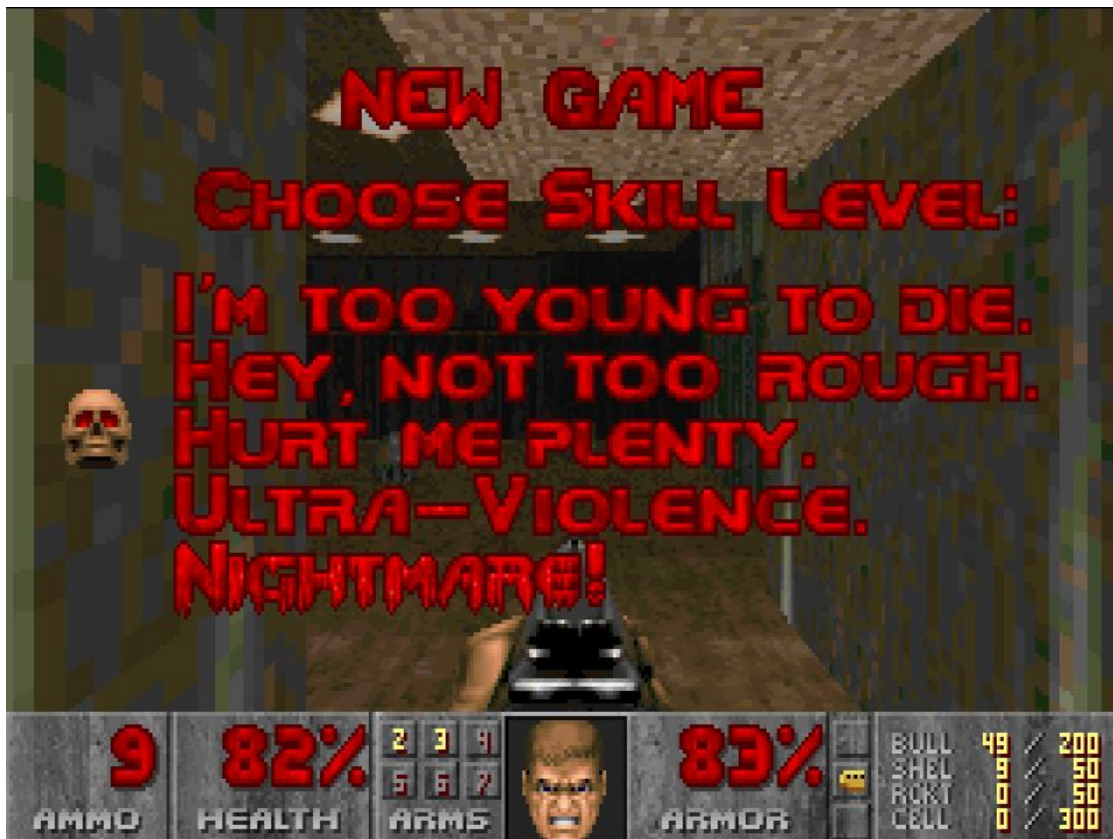
What about UDP?

Does it also work for UDP? 🤔

```
system("sysctl -w net.ipv4.ip_local_port_range='60000 60009'")  
  
s = socket(AF_INET, SOCK_DGRAM)  
s.setsockopt(SOL_IP, IP_BIND_ADDRESS_NO_PORT, 1)  
s.setsockopt(SOL_IP, IP_LOCAL_PORT_RANGE, [60000, 60009])  
s.bind(("192.0.2.1", 0)) # let the kernel find a free 4-tuple  
s.connect(("1.1.1.1", 53))
```

What about UDP?

NOPE



No source port sharing with late autobind

`IP_BIND_ADDRESS_NO_PORT` works
but not as you expect it

```
# if there is just one ephemeral port
system("sysctl -w net.ipv4.ip_local_port_range='60000 60000'")
```

```
s1 = socket(AF_INET, SOCK_DGRAM)
s1.setsockopt(SOL_IP, IP_BIND_ADDRESS_NO_PORT, 1)
s1.bind(("192.0.2.1", 0))
s1.connect(("1.1.1.1", 53))
```

```
s2 = socket(AF_INET, SOCK_DGRAM)
s2.setsockopt(SOL_IP, IP_BIND_ADDRESS_NO_PORT, 1)
s2.bind(("192.0.2.1", 0))
s2.connect(("1.0.0.1", 53)) # ⇒ EAGAIN
```

Can we make it work?

```
int inet_dgram_connect(struct socket *sock, struct sockaddr *uaddr,
                      int addr_len, int flags)
{
    // ...
    if (data_race(!inet_sk(sk)->inet_num) && inet_autobind(sk))
        return -EAGAIN;
    return sk->sk_prot->connect(sk, uaddr, addr_len);
}
```



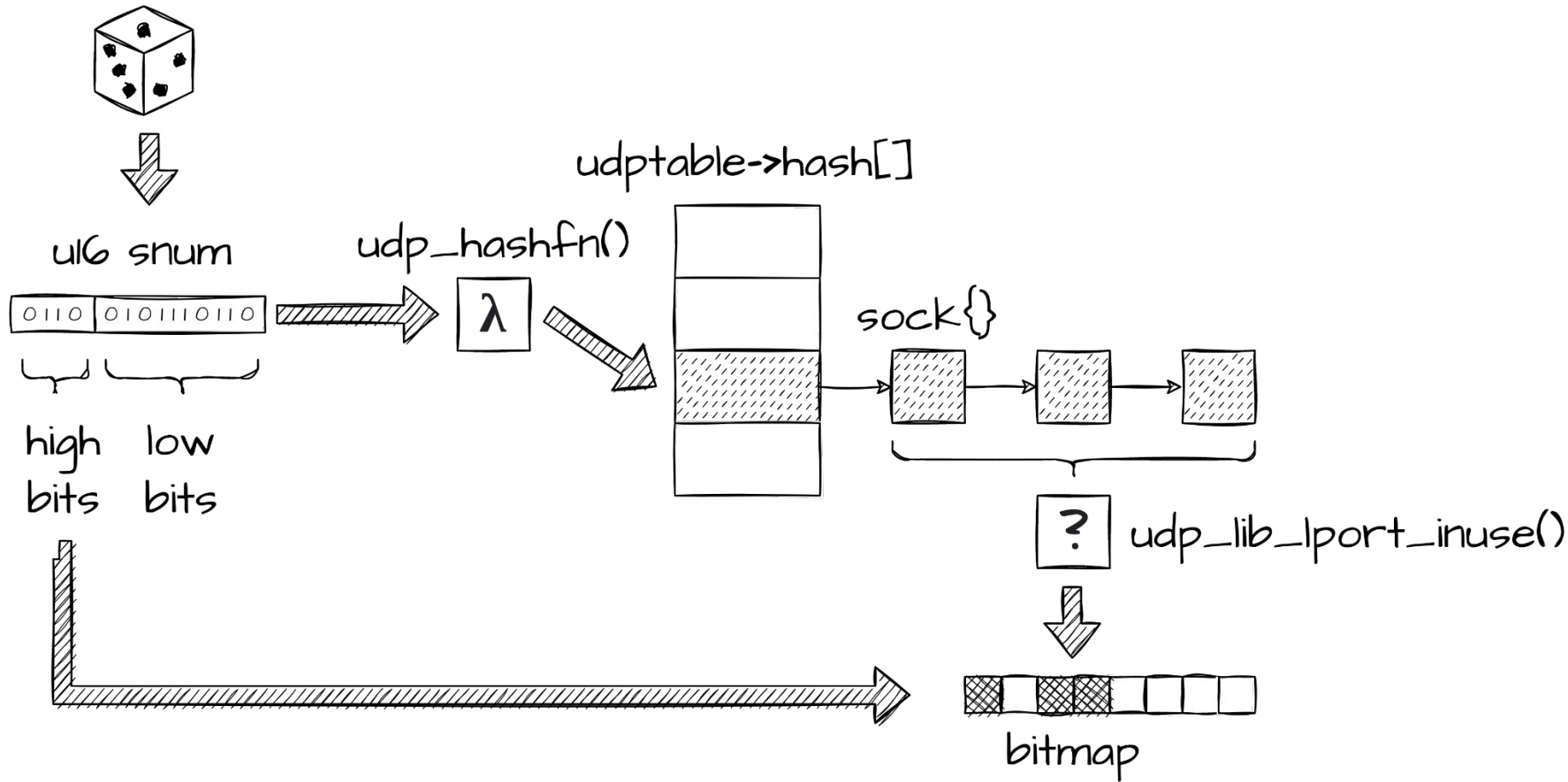
`sk->sk_prot->get_port(sk, snum)`

Can we make it work?

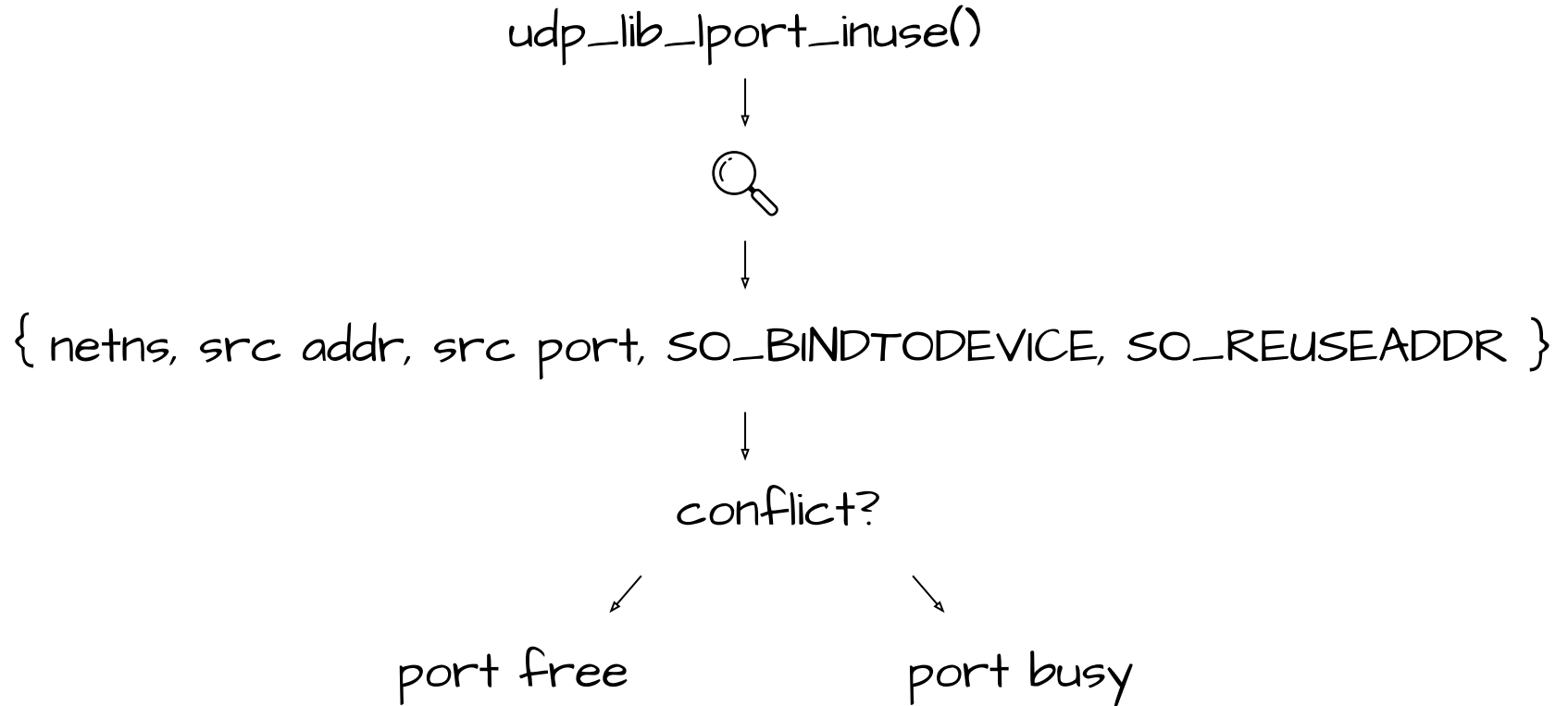
`sk->sk_prot->get_port(sk, snum)`

```
int udp_lib_get_port(struct sock *sk, unsigned short snum,
                    unsigned int hash2_nulladdr)
{
    if (!snum) {
        // (A) port not specified - find a free one
    } else {
        // (B) port given - check if free
    }
found:
    // publish the new binding
}
```

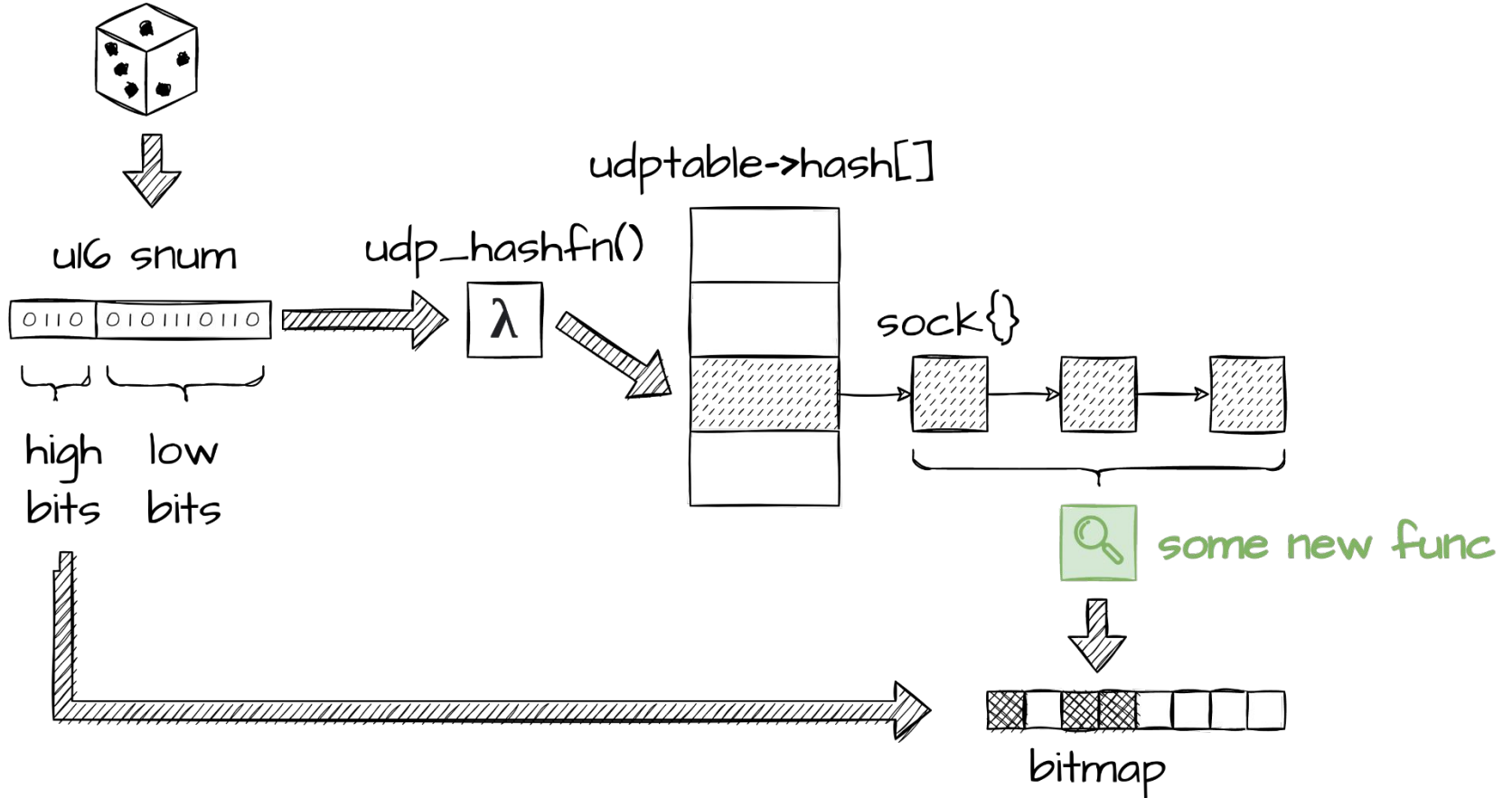
UDP get_port algorithm



UDP get_port algorithm



How can we detect 4-tuple conflicts?



How can we detect 4-tuple conflicts?

some_new_func()



{ netns, src addr, src port, SO_BINDTODEVICE, SO_REUSEADDR,
dst addr, dst port }



conflict?



port free



port busy

We could do it like so...

```
int (*bind_add)(struct sock *sk, struct sockaddr *addr, int addr_len);
```

```
@@ -2937,6 +3028,7 @@ struct proto udp_prot = {
    .sendmsg          = udp_sendmsg,
    .recvmsg         = udp_recvmsg,
    .sendpage        = udp_sendpage,
+   .bind_add         = udp_v4_bind_add,
    .release_cb      = ip4_datagram_release_cb,
    .hash            = udp_lib_hash,
    .unhash          = udp_lib_unhash,
```


We could do it like so...

```
int inet_dgram_connect(struct socket *sock, struct sockaddr *uaddr,
                      int addr_len, int flags)
{
    // ...
    if (data_race(!inet_sk(sk)->inet_num) && inet_autobind(sk))
        return -EAGAIN;
    return sk->sk_prot->connect(sk, uaddr, addr_len);
}
```

We could do it like so...

```
int inet_dgram_connect(struct socket *sock, struct sockaddr *uaddr,
                      int addr_len, int flags)
{
    // ...
    if (data_race(!inet_sk(sk)->inet_num)) { // late bind
        if (some condition) {
            if (sk->sk_prot->bind_add(sk, uaddr, addr_len))
                return -EAGAIN;
        } else {
            if (inet_autobind(sk)) // sk_prot->get_port()
                return -EAGAIN;
        }
    }
    return sk->sk_prot->connect(sk, uaddr, addr_len);
}
```

We could do it like so...

```
if (some condition) {  
    if (sk->sk_prot->bind_add(sk, uaddr, addr_len))  
        return -EAGAIN;  
}
```

```
some condition = source IP != wildcard &&  
                SO_REUSEADDR &&  
                IP_BIND_ADDRESS_NO_PORT
```

User API change

IF 4-tuple is not unique AND we are out of ephemeral ports THEN

SO_REUSEADDR	IP_BIND_ADDRESS_NO_PORT	socket(AF_INET, SOCK_DGRAM) bind(192.0.2.1, 0) connect(1.1.1.1, 53)
F	F	bind ⇒ EADDRINUSE
F	T	bind ⇒ OK connect ⇒ EAGAIN
T	F	bind ⇒ OK connect ⇒ OK conflict!
T	T	bind ⇒ OK connect ⇒ OK EAGAIN

Part V

The Alternatives

```
s = socket(AF_INET, SOCK_DGRAM)
s.setsockopt(SOL_SOCKET, SO_REUSEADDR, 1) # I'm a multicast socket!
s.setsockopt(SOL_UDP, UDP_UNICAST, 1)    # No, I'm a unicast socket!
s.bind(('127.1.1.1', 1111))
s.connect(('127.2.2.2', 2222)) # EADDRINUSE if 4-tuple conflict
```

connectx(2)

BSD System Calls Manual

connectx(2)

NAME

connectx -- initiate a connection on a socket

SYNOPSIS

```
#include <sys/socket.h>
```

```
int
connectx(int socket, const sa_endpoints_t *endpoints,
          sae_associd_t associd, unsigned int flags, const struct iovec *iov,
          unsigned int iovcnt, size_t *len, sae_connid_t *connid);
```

DESCRIPTION

The parameter *socket* is a socket. In general, **connectx()** may be used as a substitute for cases when **bind(2)** and **connect(2)** are issued in succession, as well as a mechanism to transmit data at connection establishment time.

Yet Another BPF Hook – post-connect()



```
struct bpf_sock_tuple tuple = {};
tuple.ipv4.saddr = ctx->user_ip4;
tuple.ipv4.sport = ctx->user_port;
tuple.ipv4.daddr = sk->src_ip4;
tuple.ipv4.dport = bpf_htons(sk->src_port);

struct bpf_sock *nsk = bpf_sk_lookup_udp(ctx, &tuple, sizeof(tuple.ipv4),
                                         BPF_F_CURRENT_NETNS, 0);

if (!nsk) {
    return 1; /* Conflict with existing socket. Fail connect(). */
}
```


Thank you

Jakub

jakub@cloudflare.com

@jkbs0

Marek

marek@cloudflare.com

@majek04