



Backpressure for Sockmap based Redirection

LSF/MM/BPF 26, Zagreb | Hemanth Malla

Azure Container Networking

Whats Sockmap being used for ?

- L7 Parsing
- Socket based Policing
- Tetragon
- Meta
- Cilium (A long time ago)



Whats Sockmap being used for ?

- L7 Parsing
- Socket based Policing
- Tetragon
- Meta
- Cilium (A long time ago)

But not many users for redirection in prod !



Motivation - Cilium



Next great local pod

Why Not ?

- TCP Backpressure completely broken
- Misc. Bugs
- Performance issues with small messages

Why Not?

- TCP Backpressure completely broken
- Misc. Bugs
- Performance issues with small messages

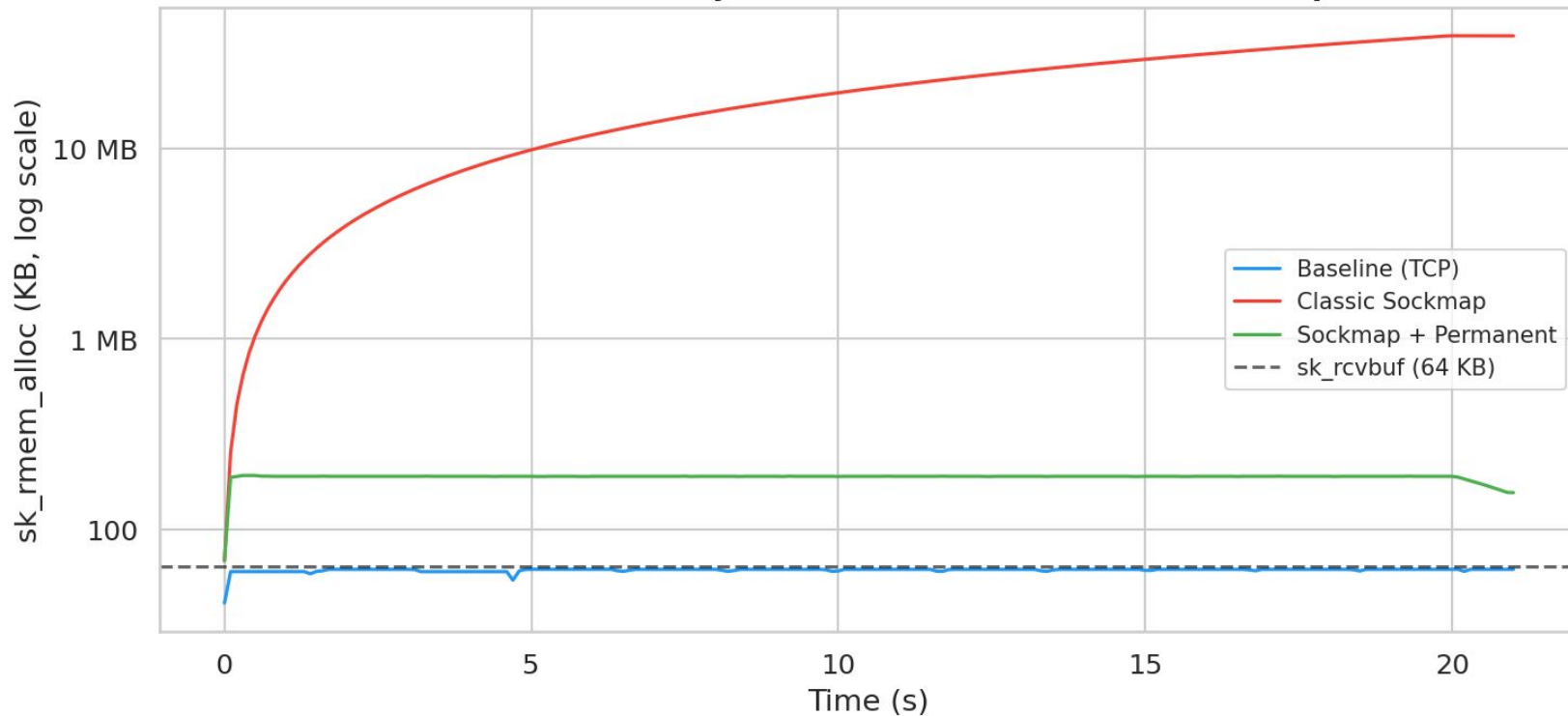


Quick Test

- Two Producers
 - Slow - 200 KB/s
 - Fast - 2000 KB/s
- Slow Consumer

Quick Tests - Unbounded Memory growth

Receiver Buffer Memory Over Time — Three Scenarios Compared



Why ?

kernel: sockmap support for blocking redirects #6438

✓ Closed



jrfastab opened on Dec 11, 2018

Member



Sockmap doesn't support blocking the rationale is that it could induce head of line blocking. However in the point to point case (or if the BPF author "knows" its ok) its useful to introduce a blocking flag to the redirect API.

Create sub-issue

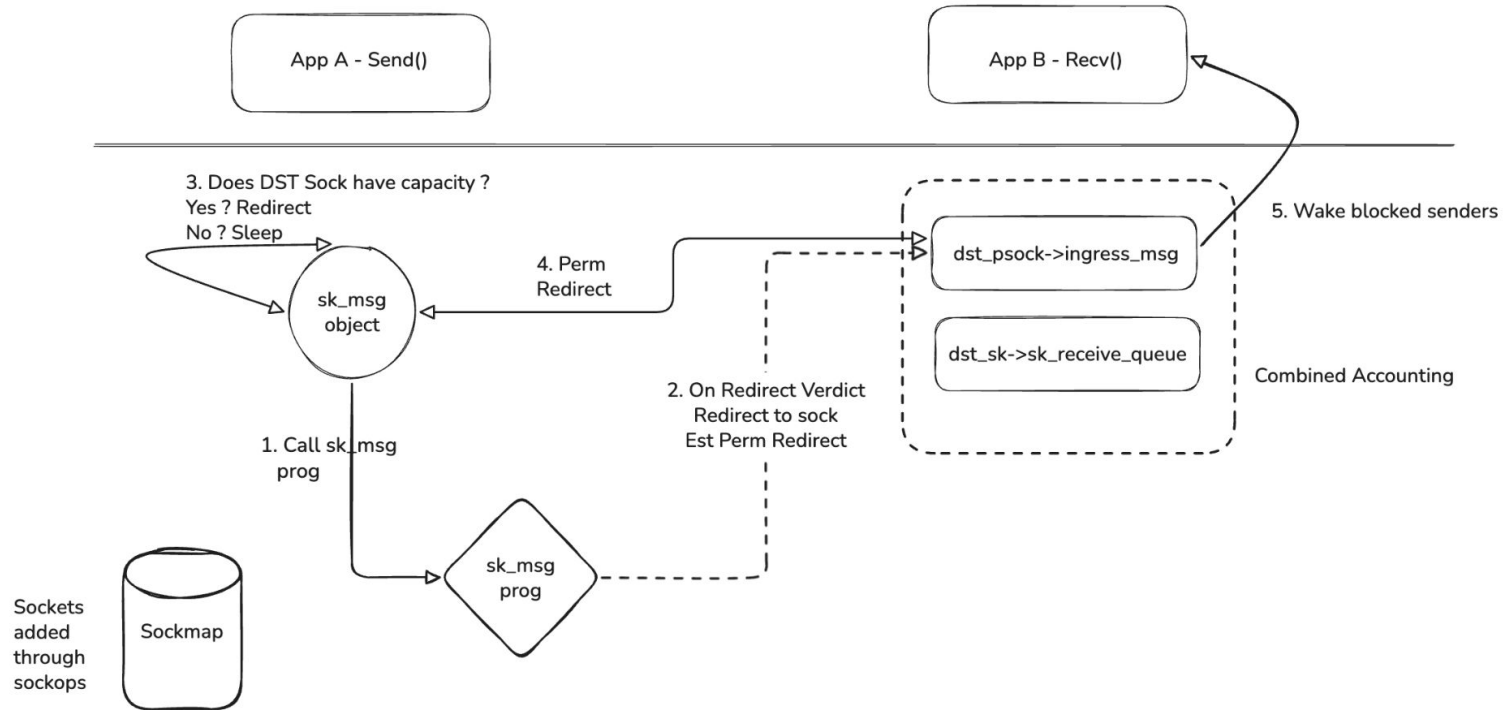


Solution - BPF_F_PERMANENT

```
SEC("sk_msg")
int redirect_prog(struct sk_msg_md *msg)
{
    int key = 1;

    return bpf_msg_redirect_hash(msg, &sock_hash, &key,
                                BPF_F_INGRESS | BPF_F_PERMANENT);
}
```

Solution - BPF_F_PERMANENT



Solution - BPF_F_PERMANENT

- The first SK_MSG verdict runs BPF normally.
- Reserves a source-to-target link if bpf prog returns perm redirect
- The first redirect must queue bytes before the link becomes active.
- Once active, future sends on the source socket bypass bpf prog and go directly to the target.
- If the target receive buffer is full -> sender waits on the target socket.
- When the target application calls recv(), receive memory is uncharged and the blocked sender is woken.

Reserve then commit ?

Should we wait till first successful data transfer ?

- Bpf_msg_redirect_map / bpf_msg_redirect_hash could be called multiple times from bpf prog
- Reserve from sk_psock_msg_verdict
- Commit from tcp_bpf_send_verdict ?



sleep + wake_up_interruptible

- Check space on dst socket in `bpf_tcp_ingress`
 - On buffer full, register on target socket's waitqueue and sleep
- Wake up blocked sender in `__sk_msg_rcvmsg()`



SK_MSG vs SK_SKB

- sk_msg in sender context
- sk_skb in receiver context
 - Sender cannot sleep
 - Sender may be remote host
 - Rely on backpressure from TCP stack on source socket buffer fill up
 - Check for dst buffer capacity before converting to sk_msg
 - On recv drain from src buffer to dst psock queue
 - Schedule backlog worker / inline ?

kTLS special handling ?

- Do we need to wait for full msg size in target buffer ?

Sockmap as the UX for redirection ?

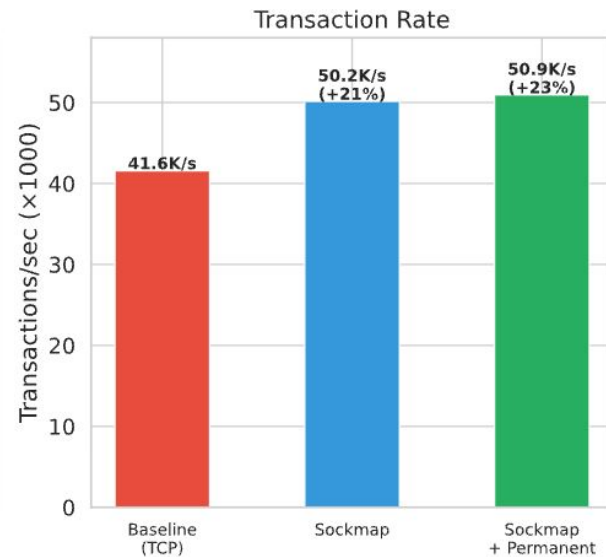
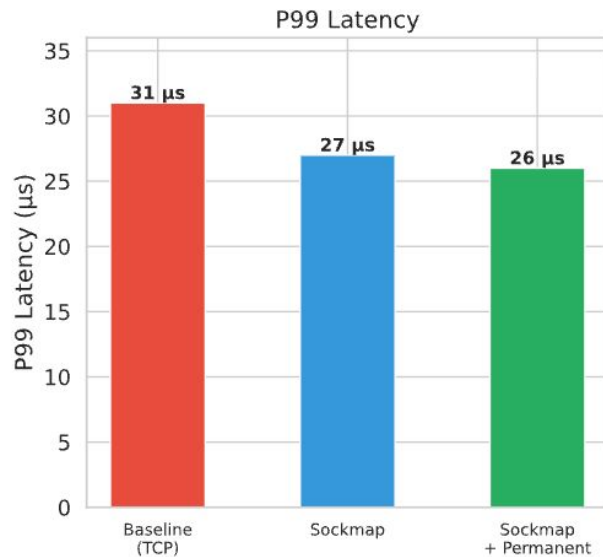
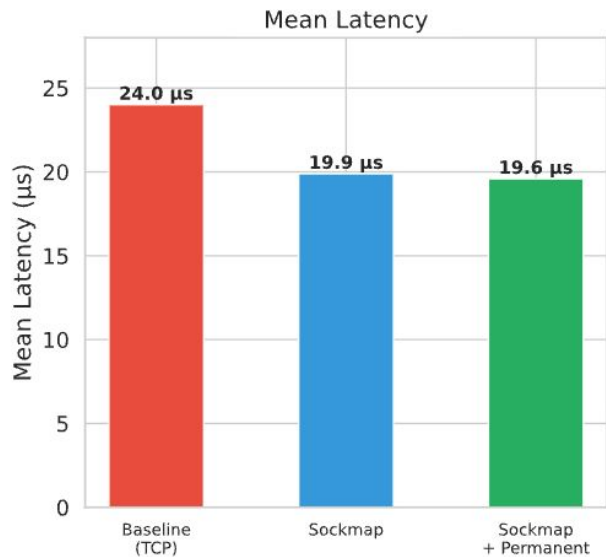
Can we decouple from the map ?

Do we care about one to many ?

Other Bugs / Blockers for adoption ?

PoC Results - TCP_RR

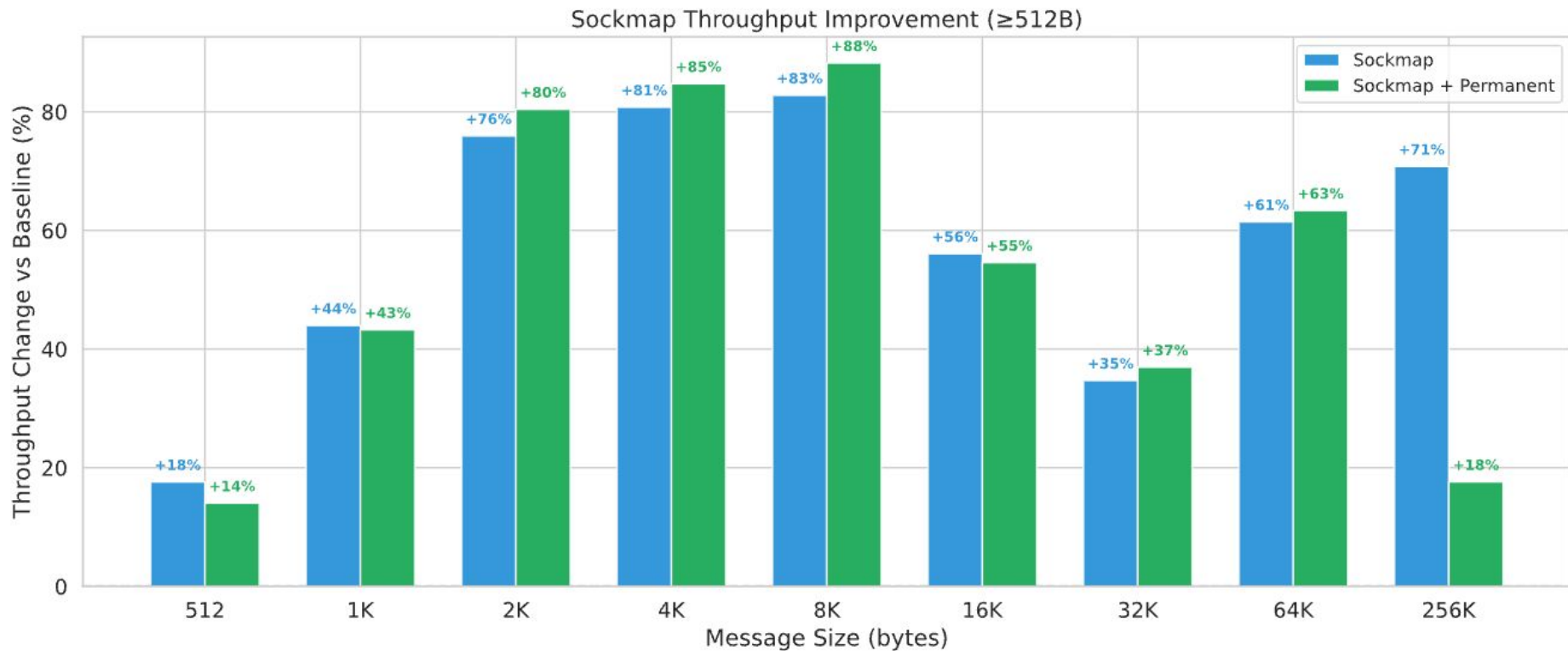
TCP_RR Summary (1B request/response)



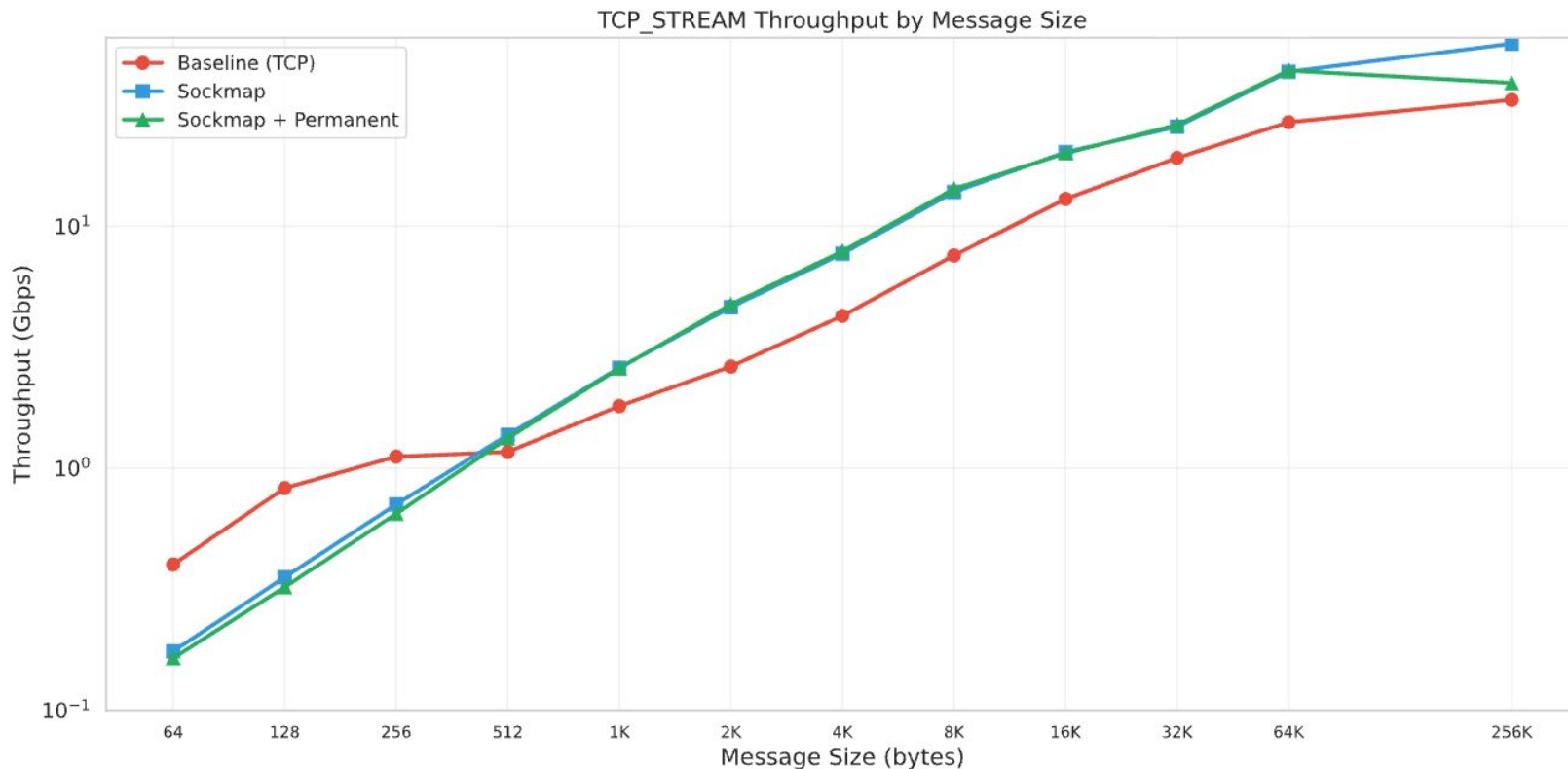
Also see old benchmarks from [Cilium 1.4](#)



PoC Results - TCP_STREAM



PoC Results - TCP_STREAM



tcp_bpf: improve ingress redirection performance with message corking

From: Cong Wang <xiyou.wangcong-AT-gmail.com>
To: netdev-AT-vger.kernel.org
Subject: [Patch bpf-next v6 0/4] tcp_bpf: improve ingress redirection performance with message corking
Date: Thu, 15 Jan 2026 11:27:33 -0800
Message-ID: <20260115192737.743857-1-xiyou.wangcong@gmail.com>
Cc: hemanthmalla-AT-gmail.com, john.fastabend-AT-gmail.com, jakub-AT-cloudflare.com, zijianzhang-AT-bytedance.com, bpf-AT-vger.kernel.org, Cong Wang <xiyou.wangcong-AT-gmail.com>
Archive-link: [Article](#)

This patchset improves skmsg ingress redirection performance by a) sophisticated batching with kworker; b) skmsg allocation caching with kmem cache.

As a result, our patches significantly outperforms the vanilla kernel in terms of throughput for almost all packet sizes. The percentage improvement in throughput ranges from 3.13% to 160.92%, with smaller packets showing the highest improvements.

For latency, it induces slightly higher latency across most packet sizes compared to the vanilla, which is also expected since this is a natural side effect of batching.

Here are the detailed benchmarks:

Throughput	64	128	256	512	1k	4k	16k	32k	64k	128k	256k
Vanilla	0.17±0.02	0.36±0.01	0.72±0.02	1.37±0.05	2.60±0.12	8.24±0.44	22.38±2.02	25.49±1.28	43.07±1.36	66.87±4.14	73.70±7.15
Patched	0.41±0.01	0.82±0.02	1.62±0.05	3.33±0.01	6.45±0.02	21.50±0.08	46.22±0.31	50.20±1.12	45.39±1.29	68.96±1.12	78.35±1.49
Percentage	141.18%	127.78%	125.00%	143.07%	148.08%	160.92%	106.52%	97.00%	5.38%	3.13%	6.32%

Latency	64	128	256	512	1k	4k	16k	32k	63k
Vanilla	5.80±4.02	5.83±3.61	5.86±4.10	5.91±4.19	5.98±4.14	6.61±4.47	8.60±2.59	10.96±5.50	15.02±6.78
Patched	6.18±3.03	6.23±4.38	6.25±4.44	6.13±4.35	6.32±4.23	6.94±4.61	8.90±5.49	11.12±6.10	14.88±6.55
Percentage	6.55%	6.87%	6.66%	3.72%	5.68%	4.99%	3.49%	1.46%	-0.93%



Thank You

