# Transparent Shared Memory Communications with eBPF

**Cong Wang, A K M Fazla Mehrab**
**System Technologies & Engineering, ByteDance**

**I‖‖I ByteDance 字节跳动**

# Overview

- Problem: The overhead of TCP communication is high for co-locating applications
- Proposal:
  - Bypass TCP stack at socket layer
  - Use shared memory as the communication channel
  - Use eBPF to maintain application transparency
- Result:
  - We observed ~12% throughput improvement for container so far
  - There is still room for improvement

ByteDance字节跳动

ByteDance 字节跳动

# Shared Memory

- Shared memory is the most effective communication we could achieve in the compute system
- Shared memory is common for IPC in OS
- For networking, RDMA technology implements this by allowing one system to share the memory of another system directly without involving the CPU
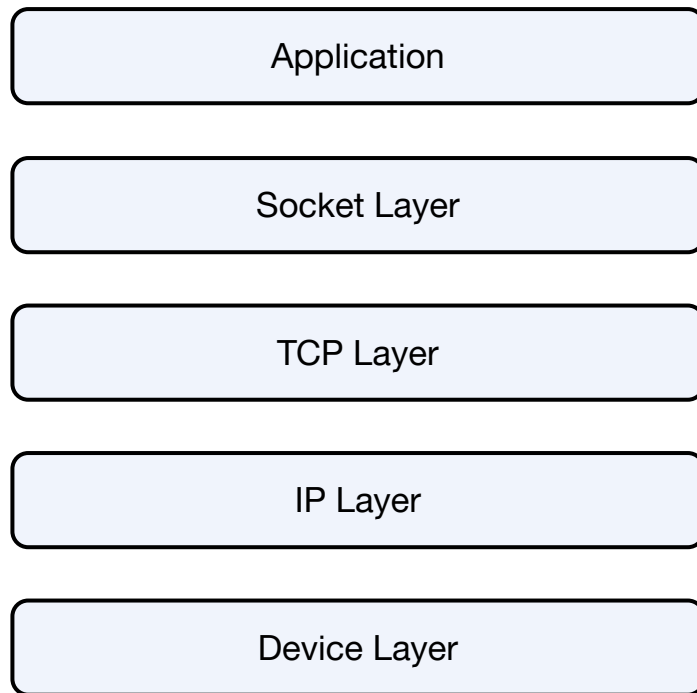- Could we leverage shared memory for TCP in a single node?

ByteDance 字节跳动

# Why still TCP?

- TCP/IP is now the de facto standard for communication in data centers and across the internet as a whole
- It is the backbone of networking in large data centers due to its reliability, scalability, and compatibility with a wide range of hardware and software systems
- TCP is the protocol of choice for applications in data centers that require reliable, connection-oriented communication over IP networks
- Its widespread adoption and support across various programming languages and platforms make it well-suited for facilitating communication between applications
- TCP socket API is the de facto standard API for networking applications

ByteDance字节跳动

ByteDance 字节跳动

# Problem Statement

Application

Socket Layer

TCP Layer

IP Layer

Device Layer

TCP stack has so many layers

ByteDance 字节跳动

# Loopback

Client Socket

Server Socket

TCP/IP

TCP/IP

Co-locating applications need to traverse stack twice

Loopback

ByteDance 字节跳动

# Co-locating Containers

# Co-resident VMs

# State-of-the-art

- Unix Domain Socket
- Virtual Socket
- RDMA
- SMC-R

ByteDance字节跳动

# Unix Domain Socket

- A method for inter-process communication (IPC) in Unix-like systems
- No stack, faster than TCP/IP for local communication
- Uses filesystem paths for socket addressing
- Supports bidirectional, streams or datagrams
- Utilized by system and user-level applications
- Operates in both connection-oriented and connectionless modes.
- Requires application modification

ByteDance字节跳动

# Virtual Socket

- Communication technology employed in virtualized and distributed environments
- Enables IPC between different VMs or containers
- Allows for reduced overhead and latency compared to traditional network-based communication
- Utilizes the host system's resources, bypassing the regular network stack
- Facilitates efficient, high-speed data transmission between distributed components
- Integral to microservices architecture and container orchestration platforms like Kubernetes
- Requires application modification

ByteDance 字节跳动

# RDMA

- **Bypasses the OS**: RDMA allows data to be transferred directly between the RAM of different computers without CPU intervention, bypassing the OS and kernel entirely
- **Low Latency and High Throughput**: Direct data transfers significantly reduce latency and increase throughput, ideal for performance-critical applications
- **Zero-Copy Networking**: Enables zero-copy networking behavior, reducing the number of data copies between applications and the network stack
- **RDMA-Capable NICs**: Requires network interface cards that support RDMA, such as those implementing InfiniBand, RoCE (RDMA over Converged Ethernet), or iWARP (Internet Wide Area RDMA Protocol)
- **Verbs API**: RDMA offers a low-level "verbs" programming API that allows for fine-grained control over RDMA operations
- **Complexity**: RDMA programming and setup can be complex and may require a deep understanding of networking concepts

ByteDance 字节跳动

# SMC-R

- **RDMA-based**: Utilizes RDMA for efficient data transfer, enabling high-speed communication between systems by bypassing the Linux kernel network stack
- **Low Latency**: Aims to reduce network latency compared to traditional TCP/IP, which is beneficial for latency-sensitive applications
- **TCP-compatibility**: Designed to be compatible with existing TCP/IP applications, allowing them to take advantage of RDMA-enabled hardware *via LD_PRELOAD*
- **Fallback to TCP**: If RDMA is not available or if setup negotiation fails, SMC-R automatically falls back to standard TCP/IP communication
- **Shared Memory**: Establishes a shared memory space between communication endpoints, allowing for efficient data exchange
- **RDMA-Capable NICs**: Requires network interfaces that support RDMA, such as those with InfiniBand or RoCE (RDMA over Converged Ethernet) capabilities
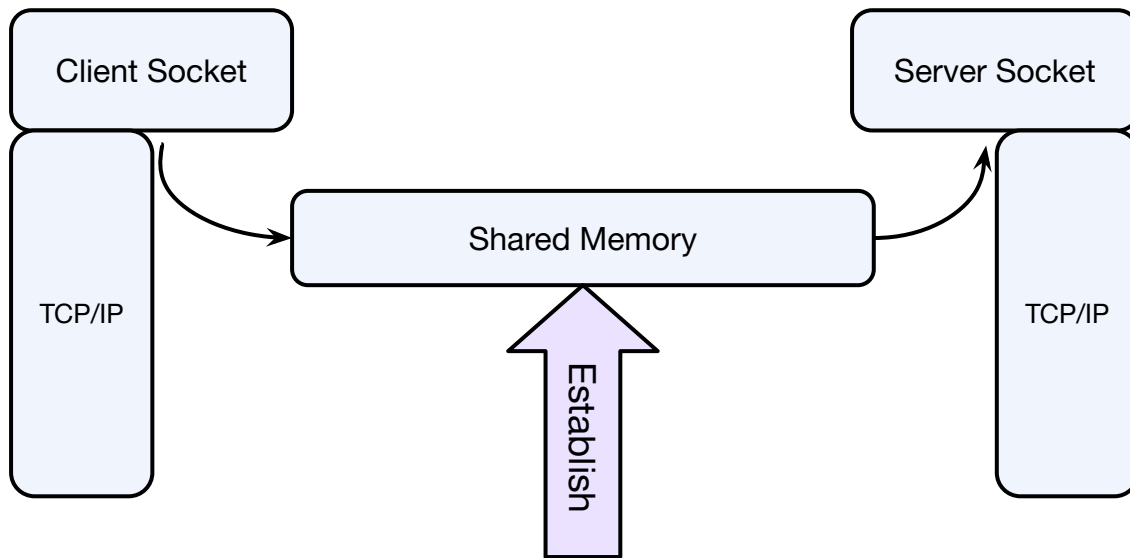
ByteDance字节跳动

ByteDance 字节跳动

# Our Idea

- Bypass TCP stack layers
  - To avoid overheads
  - With shared memory
  - Ideally zero-copy
- Maintain application transparency
  - To avoid specializations
  - To support all existing applications
- Inspirations from SMC
  - Exchange information during TCP 3-way handshake
- Inspirations from eBPF
  - Sockops
  - Sockmap

ByteDance字节跳动

# Our Idea

ByteDance字节跳动

# Non-VM Case

- Breaking it down into pieces:
  - Hijacking 3-way handshake: sockops
  - Communication channel: sockmap (sk_msg)
  - sendmsg() hook: BPF_SK_MSG_VERDICT for redirection
- It turns out Cilium already has a similar implementation: sockops-enable option
- Surprisingly, its performance is ***much worse*** than TCP!!

ByteDance字节跳动

# Cilium Socket Acceleration

- TCP is not as bad as it appears
- Linux TCP/IP stack has been optimized for decades, batching is excellent
- sk_msg is not optimized at all, not as sophisticated as skb at batching
  - For example, batching in release_sock()
- Sender needs to acquire receiver's sock lock for accounting purpose
- Sock lock becomes the source of all evil
- We have an idea for optimization and Zijian Zhang already finished preliminary work

ByteDance字节跳动

# Optimizing sk_msg

Sender

Receiver

```
lock_sock()
move skmsg to destination
Wake up sleeper
unlock_sock()
```

```
(woken up)
lock_sock()
check skmsg queue
read skmsg
unlock_sock()
(sleep)
```

```
lock_sock()
move skmsg to destination
Wake up sleeper
unlock_sock()
```

```
(woken up)
lock_sock()
check skmsg queue
read skmsg
unlock_sock()
(sleep)
```

.

.

.

ByteDance字节跳动

# Optimizing sk_msg

**Sender**

| |
|---|
| queue skmsg<br>schedule worker |

| |
|---|
| queue skmsg<br>schedule worker |

.

.

.

| |
|---|
| queue skmsg<br>schedule worker |

**Worker**

| |
|---|
| (woken up)<br>lock_sock()<br>move skmsg to destination<br>unlock_sock()<br>Wake up sleeper<br>(sleep) |

| |
|---|
| (woken up)<br>lock_sock()<br>move skmsg to destination<br>unlock_sock()<br>Wake up sleeper<br>(sleep) |

**Receiver**

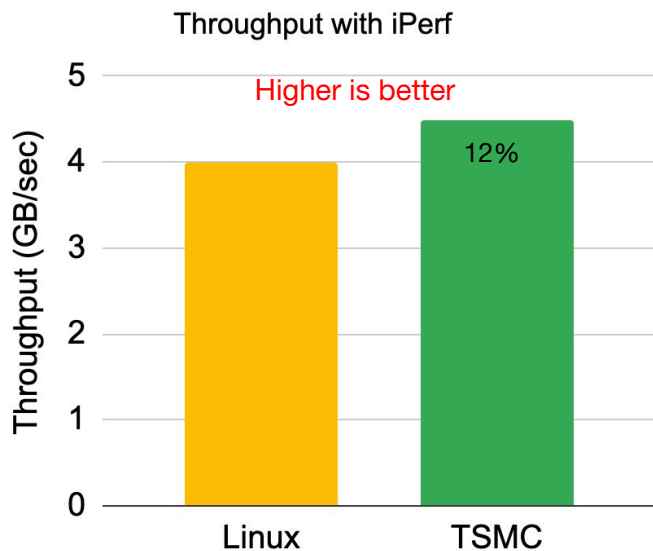| |
|---|
| (woken up)<br>lock_sock()<br>check skmsg queue<br>read skmsg<br>unlock_sock()<br>(sleep) |

| |
|---|
| (woken up)<br>lock_sock()<br>check skmsg queue<br>read skmsg<br>unlock_sock()<br>(sleep) |

ByteDance 字节跳动

# Evaluation



Throughput with iPerf — Higher is better — Linux 4, TSMC 12%

Latency with sockperf — Lower is better — Linux 29%, TSMC

TSMC→Transparent Shared Memory Communications
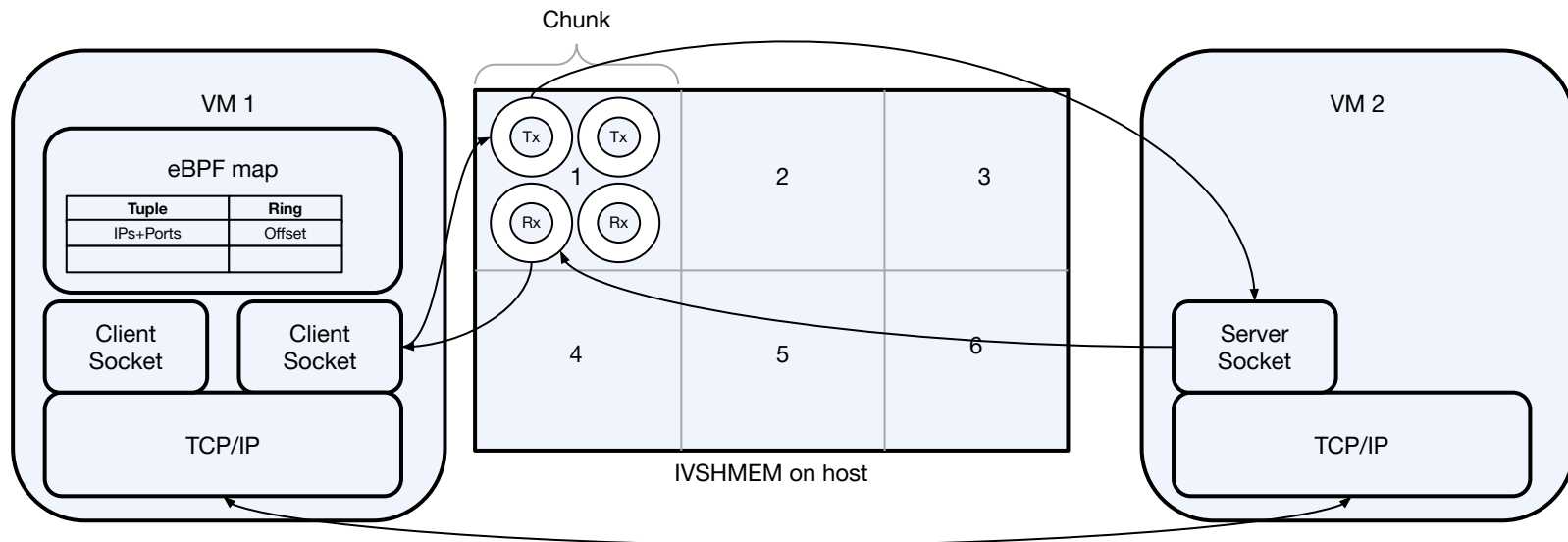
ByteDance字节跳动

# VM Case

- Breaking it down into pieces:
  - Hijacking 3-way handshake: sockops with TCP option
  - Communication channel: IVSHMEM
  - sendmsg() and recvmsg() hooks: struct_ops for struct proto
- Security assumptions:
  - VMs run known workloads in our datacenter
  - Therefore, we trust all VMs running on the same host
- No prototype implementation available yet

ByteDance字节跳动

# VM case: overview



Chunk

VM 1

eBPF map

| Tuple | Ring |
|-------|------|
| IPs+Ports | Offset |
| | |

Client Socket

Client Socket

TCP/IP

Tx  Tx

1

2

3

Rx  Rx

4

5

6

IVSHMEM on host

VM 2

Server Socket

TCP/IP

VM create
- ● Host IP
- ● Chunk offset and size

3-way handshake
- ● Host IP
- ● Tx and Rx offset

ByteDance字节跳动

# Sockops

- eBPF programs for handshake
  - Hook during SYN and ACK
  - Create Tx and Rx rings
  - Insert a new TCP option for discovery: Host IP and Ring info
  - Clean up rings if not on the same host

ByteDance 字节跳动

# IVSHMEM

- IVSHMEM
  - Creates a shared memory segment
  - Exposes it to multiple VMs as PCI device
  - VMs can map into their address spaces
- An agent on hypervisor
  - Uses IVSHMEM
  - Divides that into chunks
  - Assigns chunks during VM initialization
  - Passes host IP to VM's
- BPF arena
  - Builds an eBPF storage on top of IVSHMEM chunks
  - Backend for ring buffers

ByteDance字节跳动

# struct_ops for struct proto

- We already have many hooks in tcp_bpf_sendmsg() for sockmap
- Introduce a new struct proto with struct_ops for more flexibilities
- Use eBPF programs to implement all TCP socket operations:
  - ->sendmsg()
  - ->recvmsg()
  - ->poll()
  - ->close()
- tcp_sendmsg_sm(), tcp_recvmsg_sm():
  - Use ring buffers for sending/receiving packets
- Build an infrastructure possibly for Homa/SMC/MPTCP too
  - Still retain TCP socket APIs

ByteDance字节跳动

# Summary

| | Hardware Dependency | Application Transparency | Inter-Container Communication | Inter-VM Communication | Remote Communication |
|---|---|---|---|---|---|
| **Unix Domain Socket** | No | No | Yes (but requires shared filesystem) | No | No |
| **Vsock** | No | No | No | Yes | No |
| **RDMA** | Yes | No | No | No | Yes |
| **SMC** | No | Yes (but requires LD_PRELOAD) | No (still regular TCP) | No (but upstream is working on it) | Yes |
| **Transparent Shared Memory Communication** | No | Yes | Yes | Yes | No |

ByteDance字节跳动

# Questions?

ByteDance 字节跳动