

# XDP and page\_pool allocator

...let's go fast!

Ilias Apalodimas  
Linaro

Lorenzo Bianconi  
Red Hat

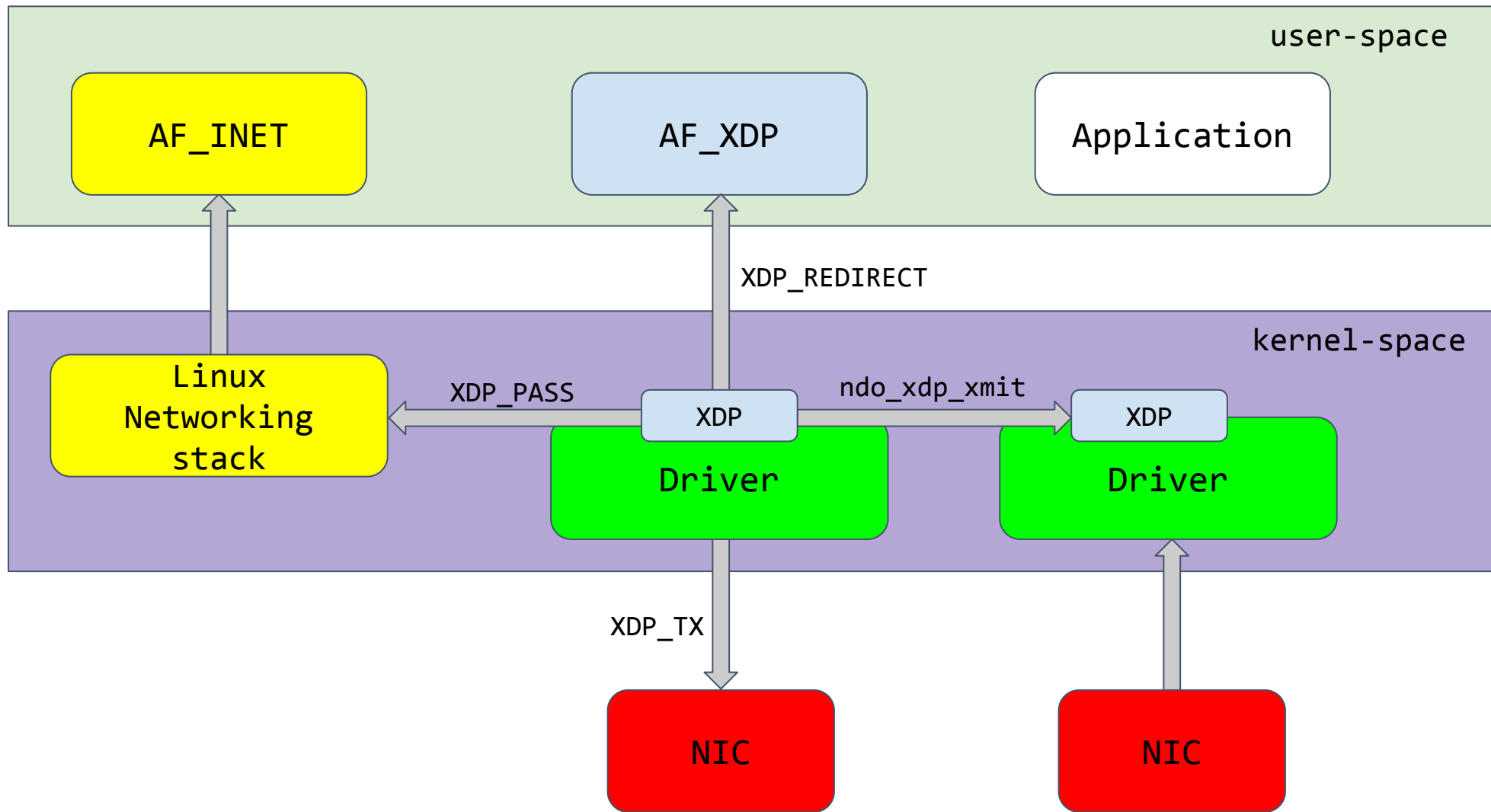
## About us:

- Ilias:
  - Serving as a co-maintainer for page\_pool API - [ilias.apalodimas@linaro.org](mailto:ilias.apalodimas@linaro.org)
  - Added XDP and page\_pool support to netsec driver
- Lorenzo:
  - Software engineer @ Red Hat, driver maintainer - [lorenzo@kernel.org](mailto:lorenzo@kernel.org)
  - Added XDP and page\_pool support to mvneta driver

# What is XDP?

XDP (eXpress Data Path) is a Linux in-kernel fast path:

- XDP can be considered as a software offload layer for the kernel networking stack
- Driver level (RX) hooks at DMA layer **before skb allocation**
  - Just after `dma_sync_*_for_cpu()`
- XDP designed for speed:
  - Operates @ L2-L3 while the networking stack works @ L2-L7
  - Skip unnecessary processing in the networking stack (route lookup, netfilter hooks, etc.)
- Not a bypass but in-kernel fast path:
  - Leverages existing kernel infrastructure
  - Programmable flexibility via eBPF sandboxing
- XDP use-cases:
  - Anti-DDoS (Facebook, CloudFlare)
  - L4 Load Balancer (Katran - Facebook)
  - etc.



# XDP requirements

XDP memory model:

- XDP frame in physical contiguous memory
  - BPF Direct-Access for validating correctness
  - Read and Write access to the DMA buffers
  - Disable jumbo frames loading a BPF program
  - No paged frames support, data cannot be split across pages
- Linear DMA pages must provide space for metadata
  - XDP headroom to push/pop header through `bpf_xdp_adjust_head()`
  - XDP\_PASS: need to reserve tailroom for `skb_shared_info` to rely on `build_skb()`
- Cannot allocate page fragments to support it (e.g. through `napi_alloc_skb()`)
- Rx buffers must be recycled to get high speed!

# page\_pool allocator

- Optimized for one packet per page
  - Supports split-page model (usually driver is in charge of recycling)
- Native buffer recycling
  - in-irq cache and ptr\_ring cache
  - Currently supported for XDP\_DROP but XDP\_PASS is coming
- Allocates order^n pages (usually order-0 = 4K page)
- Fast, usually runs in NAPI context, no extra locking overhead
- DMA management can be done via the page\_pool API
  - DMA-mapping capability (keeps page mapped)
  - DMA-sync for cache non-coherent devices

# Caveats

- Although page\_pool is faster due to native recycling for XDP, it's slower for skb
- Bigger memory footprint (linear to number of descriptors used), unless page splitting is implemented
- Working on native skb recycling, which should eliminate the skb use case penalty
  - The current mvneta driver was already allocating a page per packet. The recycling patches boost performance by ~25% on < 512b packets

“...talk is cheap. Show me the code”



- Intel and Mellanox XDP implementations are complex (naturally since it's a complex hardware)
- mvneta (marvell 1Gbit) and netsec (socionext 1Gbit) can serve as a simplified guideline on how to add XDP support
- We need all XDP verdicts covered to accept a driver
  - XDP\_DROP
  - XDP\_TX
  - XDP\_PASS
  - XDP\_REDIRECT, `ndo_xdp_xmit()`

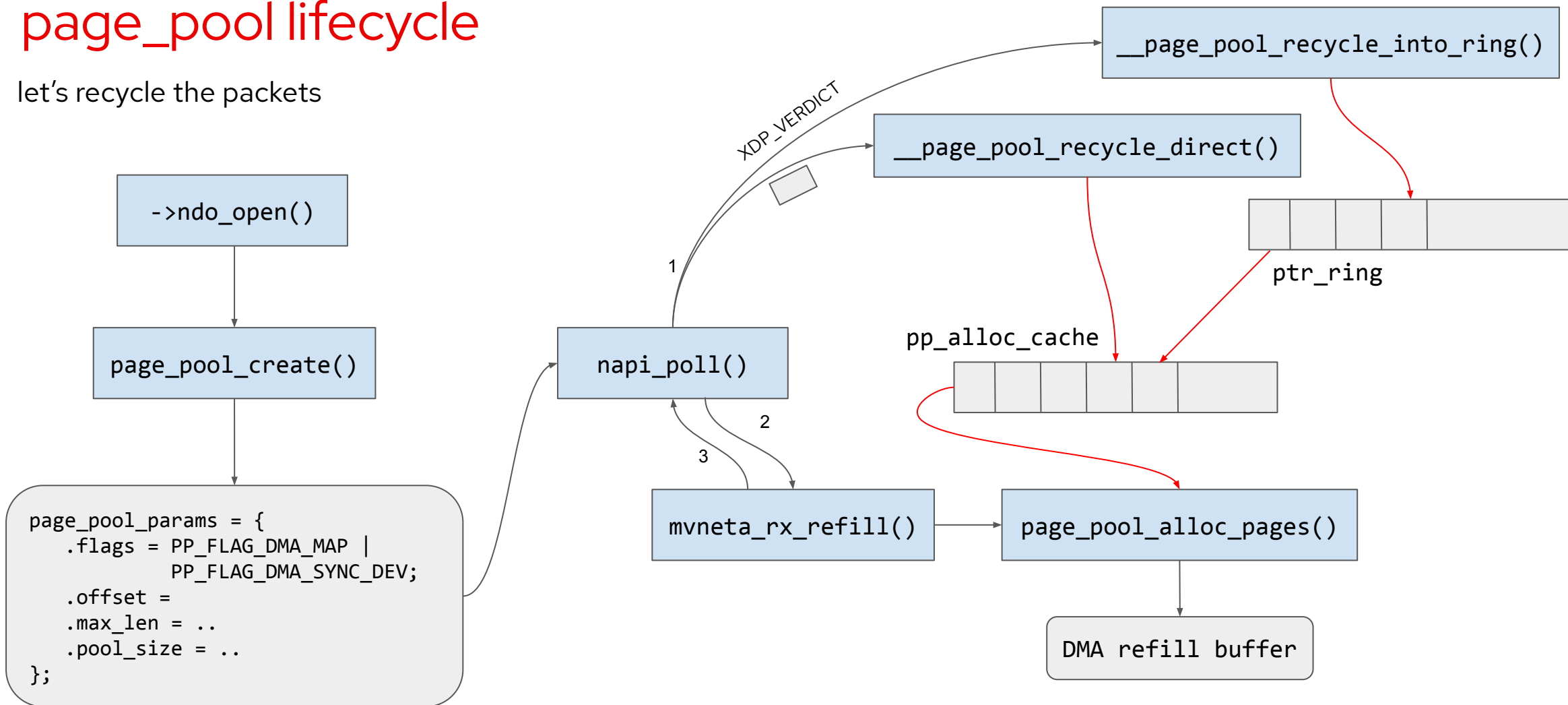


# Marvell ESPRESSObin - mvneta

SoC	Marvell Armada 3700LP (88F3720) dual core ARM Cortex A53 processor up to 1.2GHz
System Memory	1 GB DDR3 or optional 2GB DDR3
Storage	1x SATA interface 1x micro SD card slot with footprint for an optional 4GB EMMC
Network Connectivity	1x Topaz Networking Switch 2x GbE Ethernet LAN 1x Ethernet WAN 1x MiniPCIe slot for Wireless/BLE peripherals
USB	1x USB 3.0 1x USB 2.0 1x micro USB port
Expansion	2x 46-pin GPIO headers for accessories and shields with I2C, GPIOs, PWM, UART, SPI, MMC, etc.
Misc	Reset button, JTAG interface
Power supply	12V DC jack or 5V via micro USB port
Power consumption	Less than 1W thermal dissipation at 1 GHz

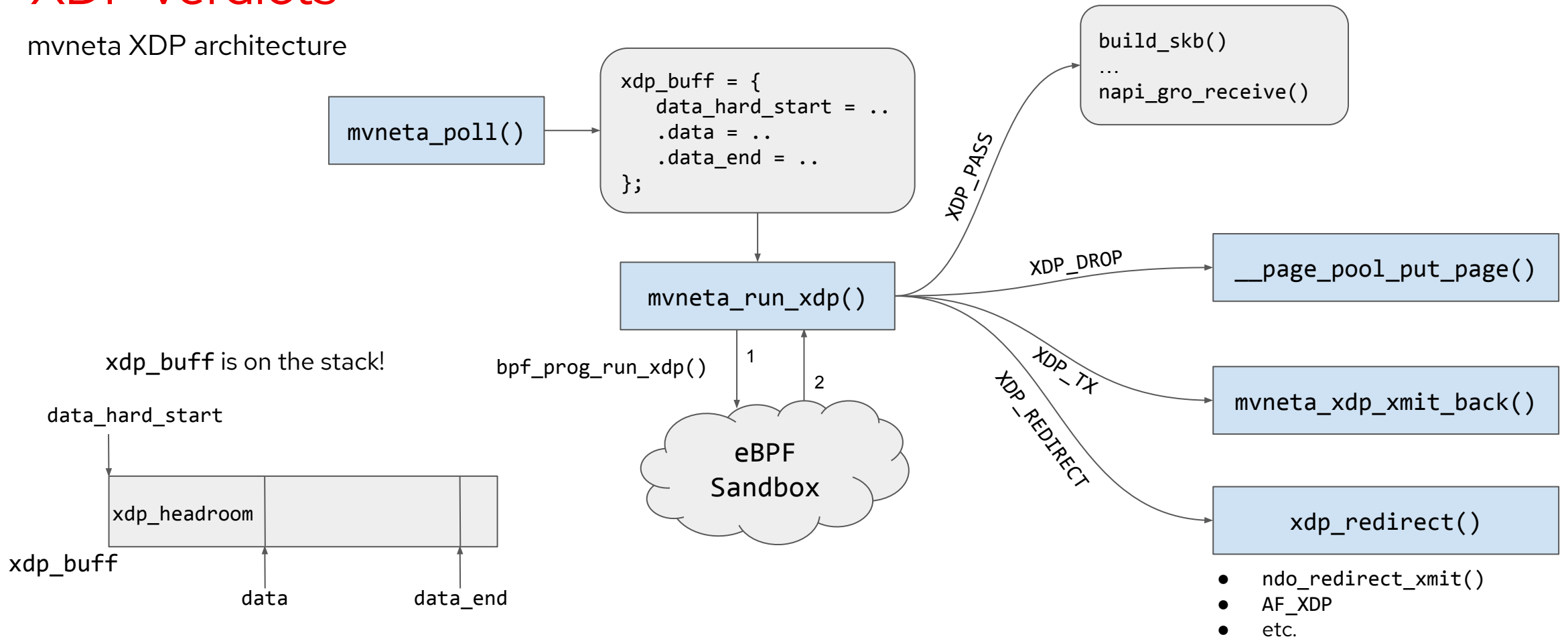
# page\_pool lifecycle

let's recycle the packets



# XDP verdicts

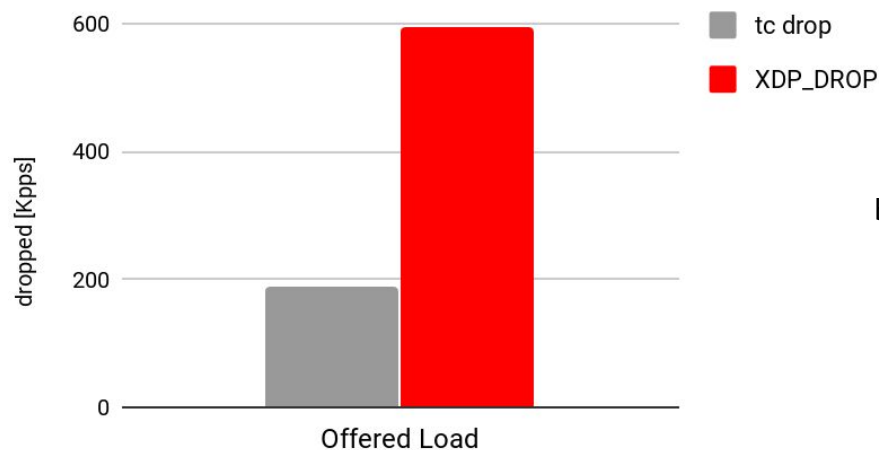
mvneta XDP architecture



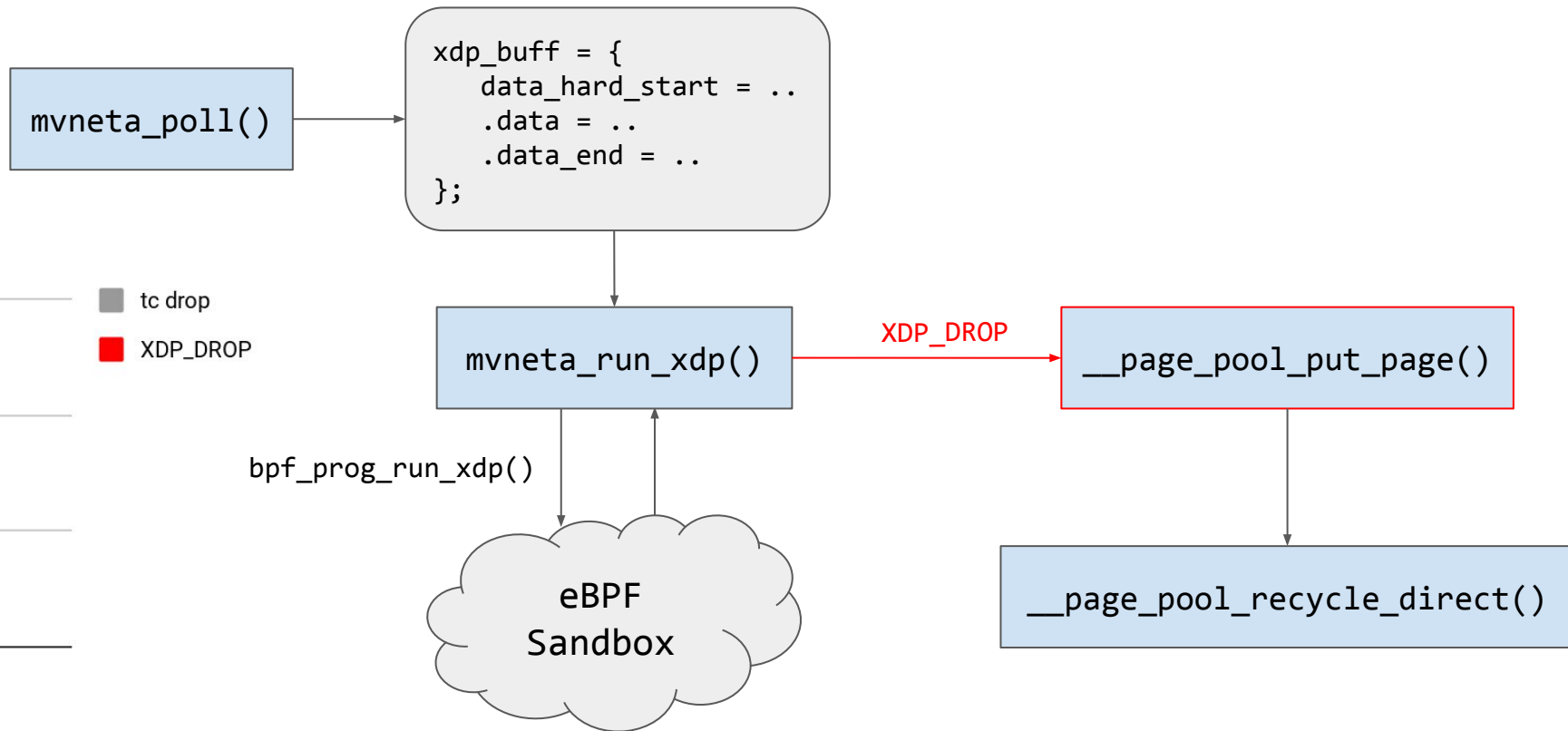
# XDP\_DROP

let's drop as fast as we can

XDP\_DROP vs tc drop

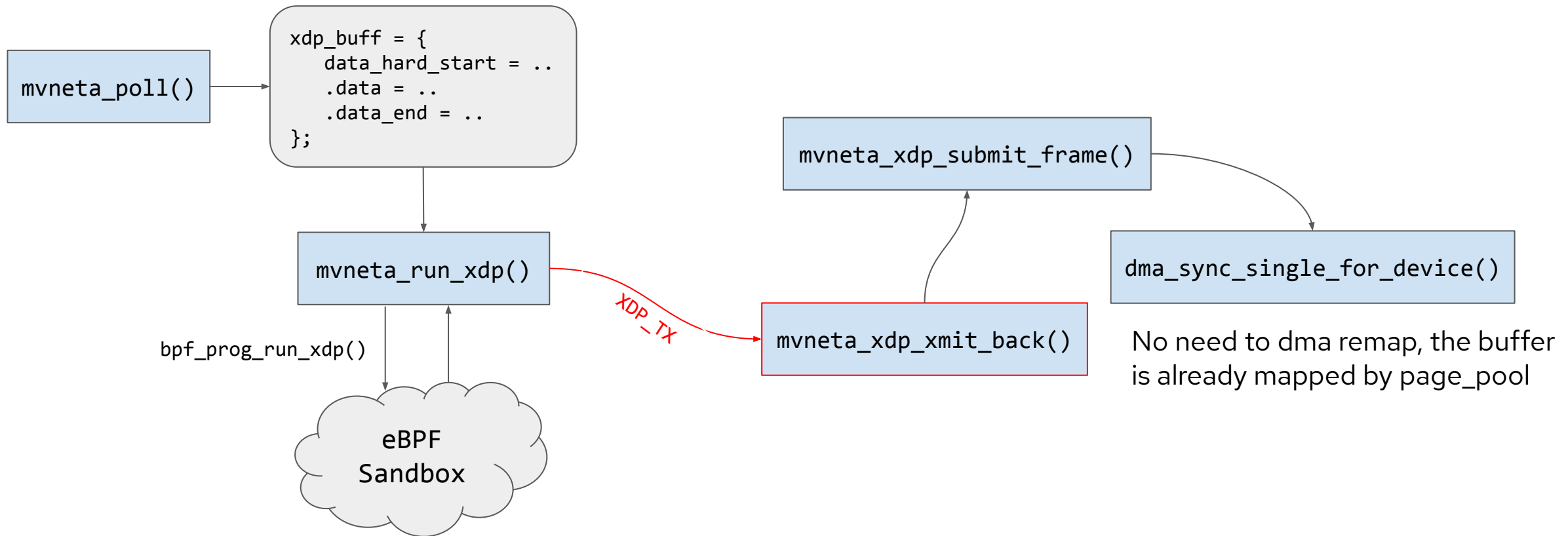


- packet size: 64B
- `$tc qdisc add dev eth0 clsact`  
`$tc filter add dev eth0 ingress matchall action gact drop`



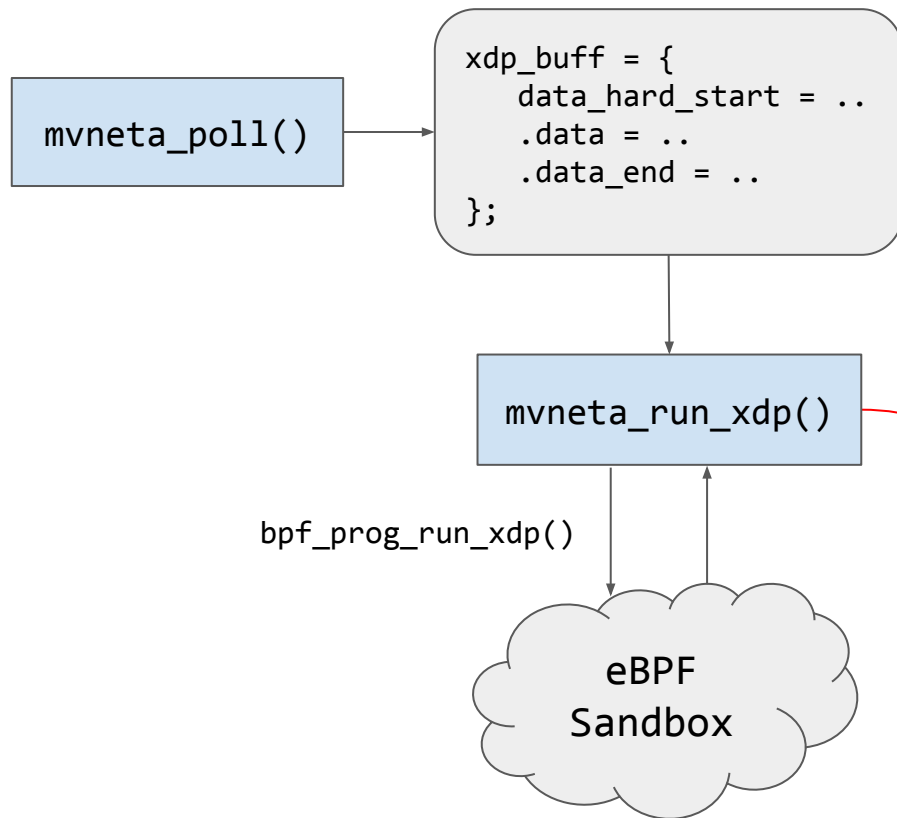
# XDP\_TX

let's send the packet back

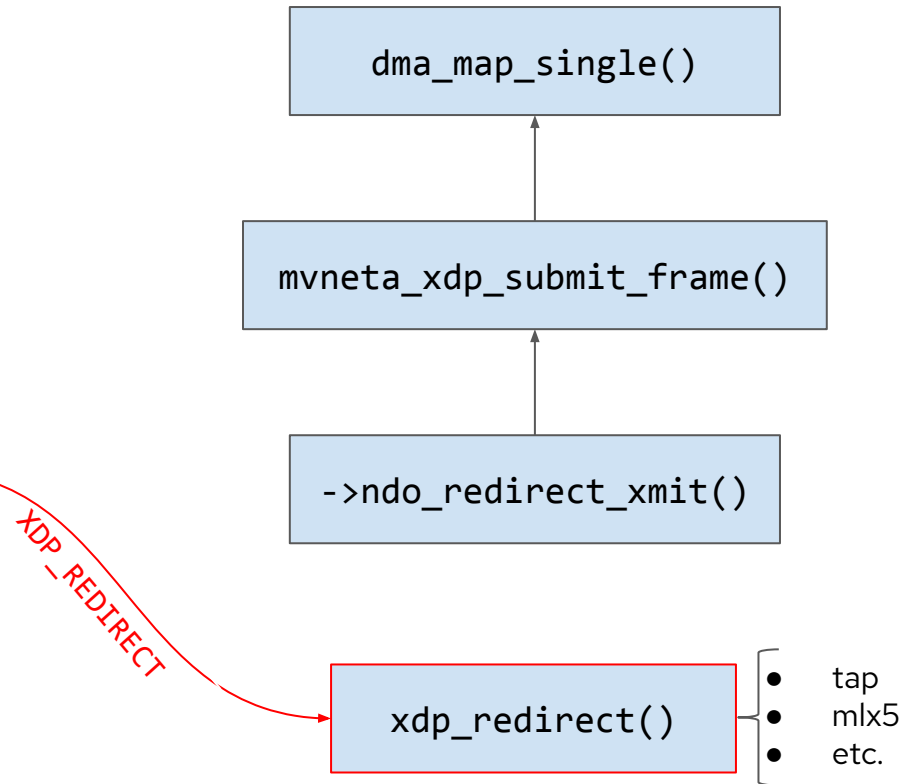


# XDP\_REDIRECT

let's forward the packet

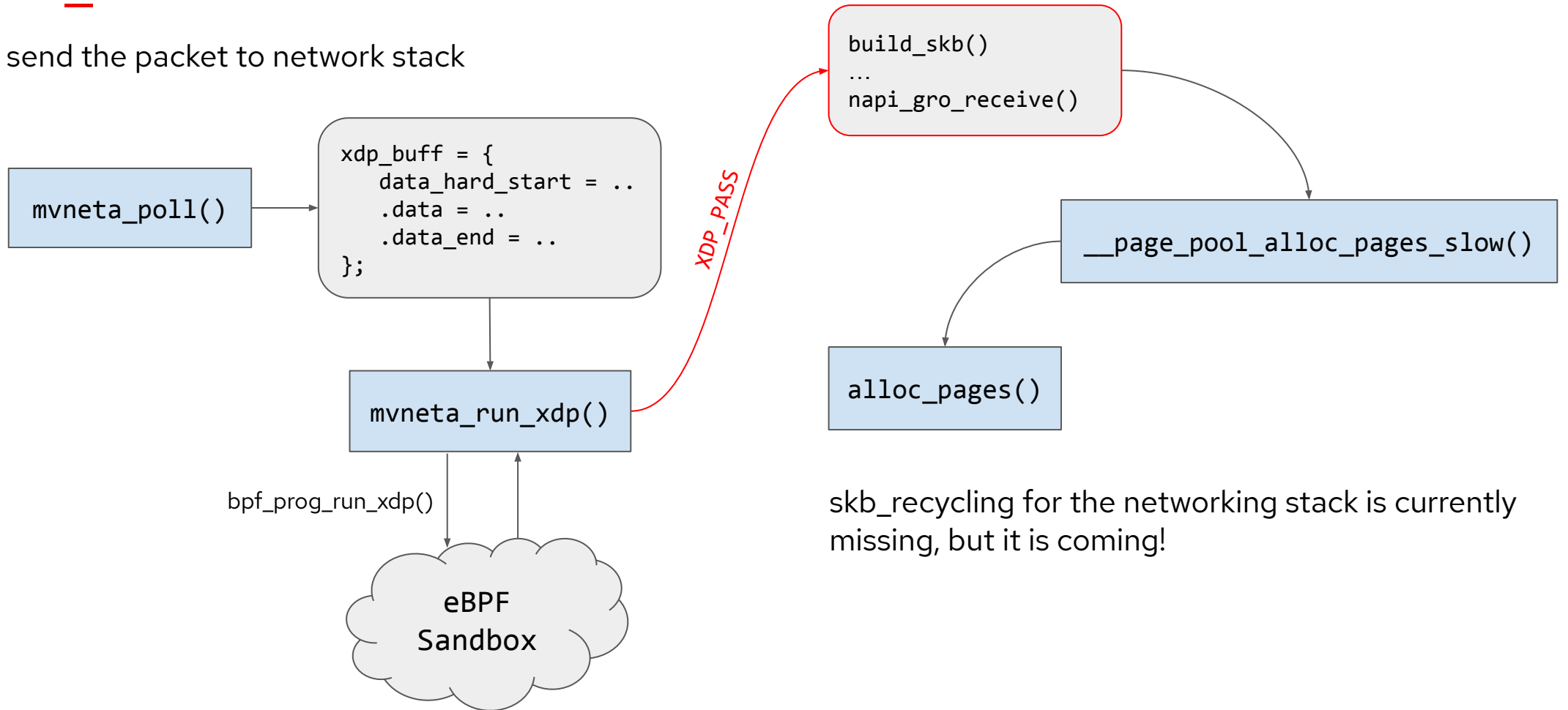


Buffer needs mapping before sending



# XDP\_PASS

let's send the packet to network stack



skb\_recycling for the networking stack is currently missing, but it is coming!

# Conclusions

- XDP can be considered as a software offload layer for the kernel networking stack
- XDP memory model
  - Contiguous memory area
- page\_pool allocator
  - DMA buffer recycling
- mvneta XDP architecture
  - XDP\_DROP
  - XDP\_TX
  - XDP\_REDIRECT
  - XDP\_PASS
- Future work
  - **skb recycling for XDP\_PASS**
  - mvneta: XDP support for hw buffer manager (e.g. ClearFog)
  - mvneta: XDP roadmap
  - mvneta: native AF\_XDP support



# Q&A:



# Thank you

- <https://github.com/xdp-project>
- Contact:
  - [ilias.apalodimas@linaro.org](mailto:ilias.apalodimas@linaro.org)
  - [lorenzo@kernel.org](mailto:lorenzo@kernel.org)