

# *Linux Networking Futures 2010*

David S. Miller

Red Hat Inc.

Boston, USA, 2010

# RPS

Linux  
Networking  
Futures 2010

David  
S. Miller

RPS and RFS

SKB  
list\_head

Iteration  
Elimination

Kill rtcache

- SMP stateless flow separation in software
- Hash computed in software
- OR taken from hardware
- Cross-calls used to doorbell remote cpu queues

# RFS

Linux  
Networking  
Futures 2010

David  
S. Miller

RPS and RFS

SKB  
list\_head

Iteration  
Elimination

Kill rtcache

- More sophisticated
- I/O calls monitored
- CPU of I/O call recorded in hash table
- Recorded CPU directs future packet steering
- Out-of-order issues avoided during hash changes

# XFS

Linux  
Networking  
Futures 2010

David  
S. Miller

RPS and RFS

SKB  
list\_head

Iteration  
Elimination

Kill rtcache

- I will leave this to Tom.

# STATE

Linux  
Networking  
Futures 2010

David  
S. Miller

RPS and RFS

SKB  
list\_head

Iteration  
Elimination

Kill rtcache

- RPS/RFS fully merged and functional
- XFS is in review state, 2.6.38 likely
- Unfortunately RPS/RFS not enabled by default
- Hey, it helps “loopback”
- Lack of infrastructure

# SKB LIST\_HEAD: HOW CLOSE?

Linux  
Networking  
Futures 2010

David  
S. Miller

RPS and RFS

SKB  
list\_head

Iteration  
Elimination

Kill rtcache

- We're almost there.
- Abstractions exist for all uses
- Even for frag\_list
- PPP, ISDN-PPP, LRO driver code, and GRO

# SKB LIST\_HEAD: NEXT STEPS

Linux  
Networking  
Futures 2010

David  
S. Miller

RPS and RFS

SKB  
list\_head

Iteration  
Elimination

Kill rtcache

- Test and merge PPP/ISDN-PPP conversion
- Move forward with ixgbe et al. frag\_list conversions
- Pull the trigger
- Fix the bugs
- Pop the champagne

# ITERATION ELIMINATION: NETIF\_RX() SUCKS

Linux  
Networking  
Futures 2010

David  
S. Miller

RPS and RFS

SKB  
list\_head

Iteration  
Elimination

Kill rtcache

- Transmit path of loopback and tunnels
- Queue packet, signal interrupt
- Expensive way to avoid stack overflow



# ITERATION ELIMINATION: QUICK HACKS

Linux  
Networking  
Futures 2010

David  
S. Miller

RPS and RFS

SKB  
list\_head

Iteration  
Elimination

Kill rtcache

- netif\_receive\_skb() with recursion limit
- Eric Dumazet
- lock\_sock sets mark
- netif\_rx() queues when mark set
- release\_sock runs packets in queue
- Issue: What is we sleep with socket lock?

# ITERATION ELIMINATION: MORE FORMAL

Linux  
Networking  
Futures 2010

David  
S. Miller

RPS and RFS

SKB  
list\_head

Iteration  
Elimination

Kill rtcache

- `net_start_transmit()`
- Acts like `lock_sock` idea above
- But never done before sleep'able section
- Finish with `net_end_transmit()`
- Runs queue etc.
- State is `cpu-local`, like current backlog

# RTCACHE ELIMINATION: WHY NOW?

Linux  
Networking  
Futures 2010

David  
S. Miller

RPS and RFS

SKB  
list\_head

Iteration  
Elimination

Kill rtcache

- Eric Dumazet's optimizations
- fib\_rie faster than routing cache
- rtcache non-determination
- rtcache complexity

# RTCACHE ELIMINATION: MAIN PROBLEM

Linux  
Networking  
Futures 2010

David  
S. Miller

RPS and RFS

SKB  
list\_head

Iteration  
Elimination

Kill rtcache

- Metrics
- Flow based metrics exist only in routing cache
- As do path based metrics (PMTU, redirects, etc.)
- Supplement with inetpeer cache?
- New metric cache by flow?
- Issue of IPSEC based metrics?