

DoS Protection in Knot Resolver

using multi-prefix query counting

Introduction: DoS

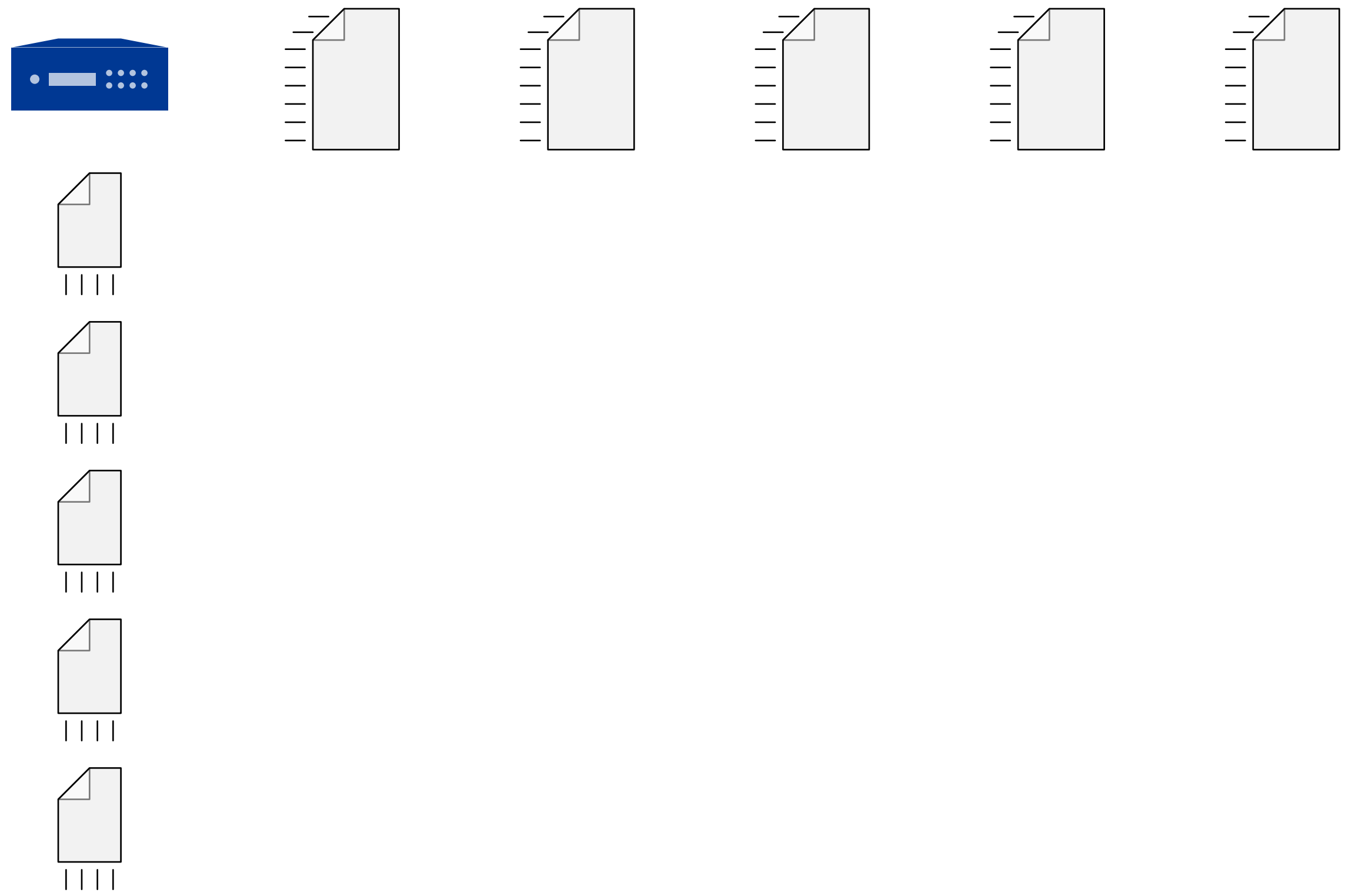
- no rate limiting ability so far, adding now in 2024
- better *inside* resolver to understand DoT, DoH, and in future DoQ
- public resolvers are the main use case, e.g.:
 - cz.nic's ODVR: <https://www.nic.cz/odvr/>
 - DNS4EU instances (to become public in 2025)
- also mostly applied in authoritative Knot DNS ≥ 3.4

Overview

- still a common DoS technique
- attacks through some UDP servers where answers are bigger than queries, therefore amplifying the attacker's traffic

1. amplification att.

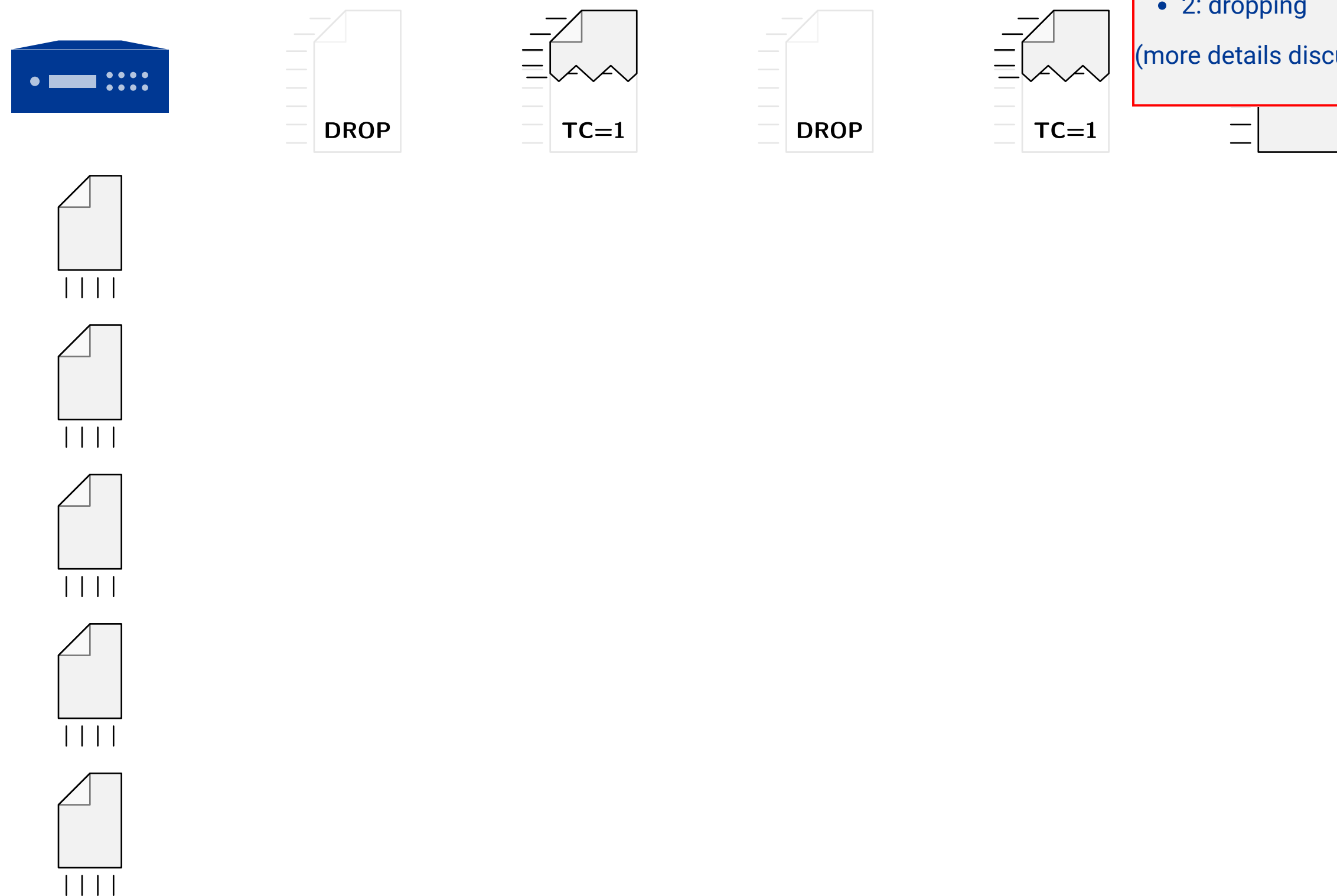
- forged UDP source
- answers >> queries



Overview

1. amplification att.

- forged UDP source
- answers >> queries
- size limit (1232 B)
- rate limiting
 - truncation
 - dropping



- Size of UDP replies was limited already, primarily to avoid issues with fragmentation.
- Now additionally: restrict response rate for each address/network, to protect *them*
- 1: truncation
 - same length of answer as its query, i.e. not amplifying
 - sane clients retry over TCP; there you can't forge source IP like that
- 2: dropping

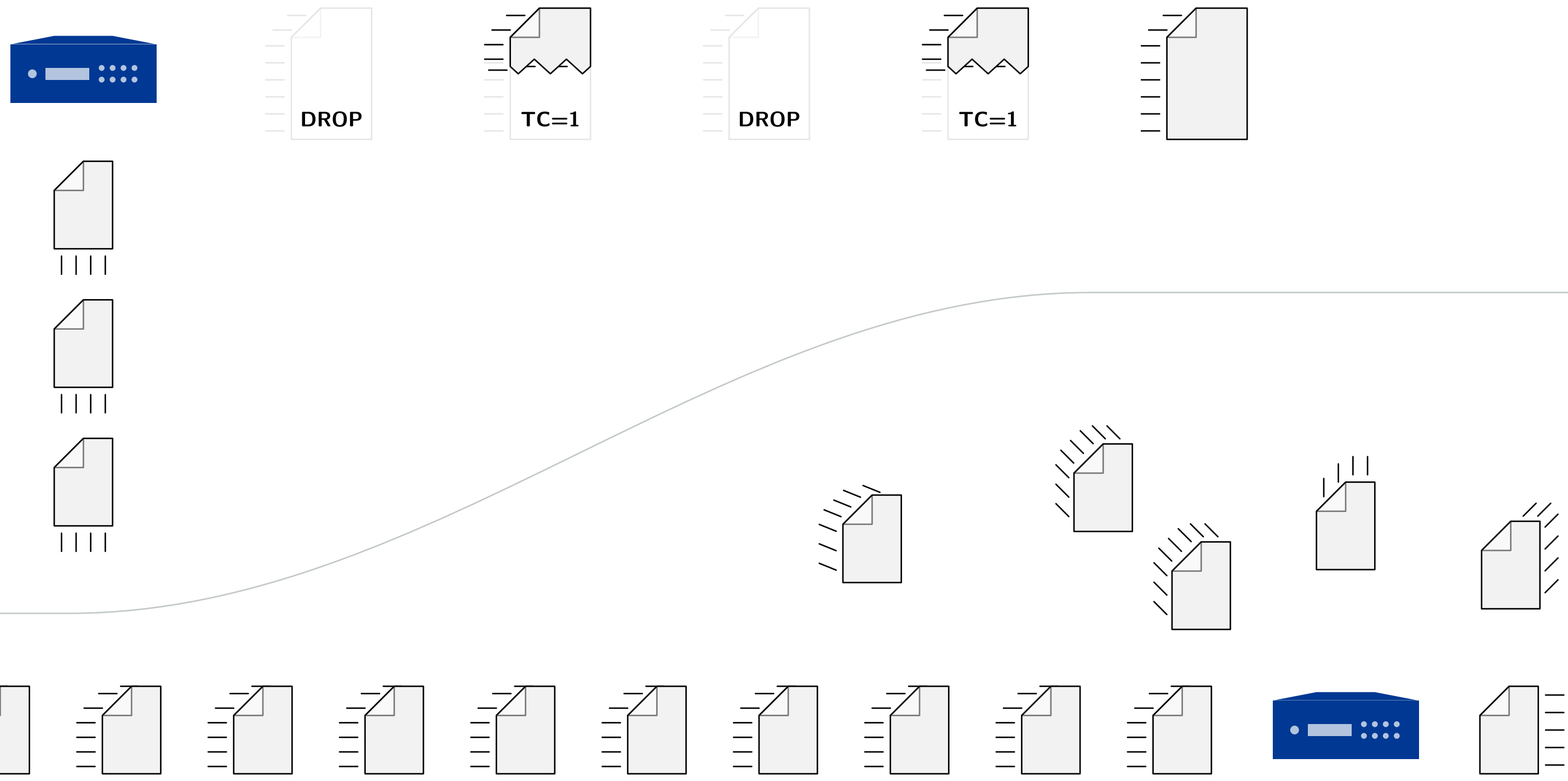
(more details discussed later)

• Many requests from the same origin can exhaust cpu time...

Overview

1. amplification att.

- forged UDP source
- answers >> queries
- size limit (1232 B)
- rate limiting
 - truncation
 - dropping

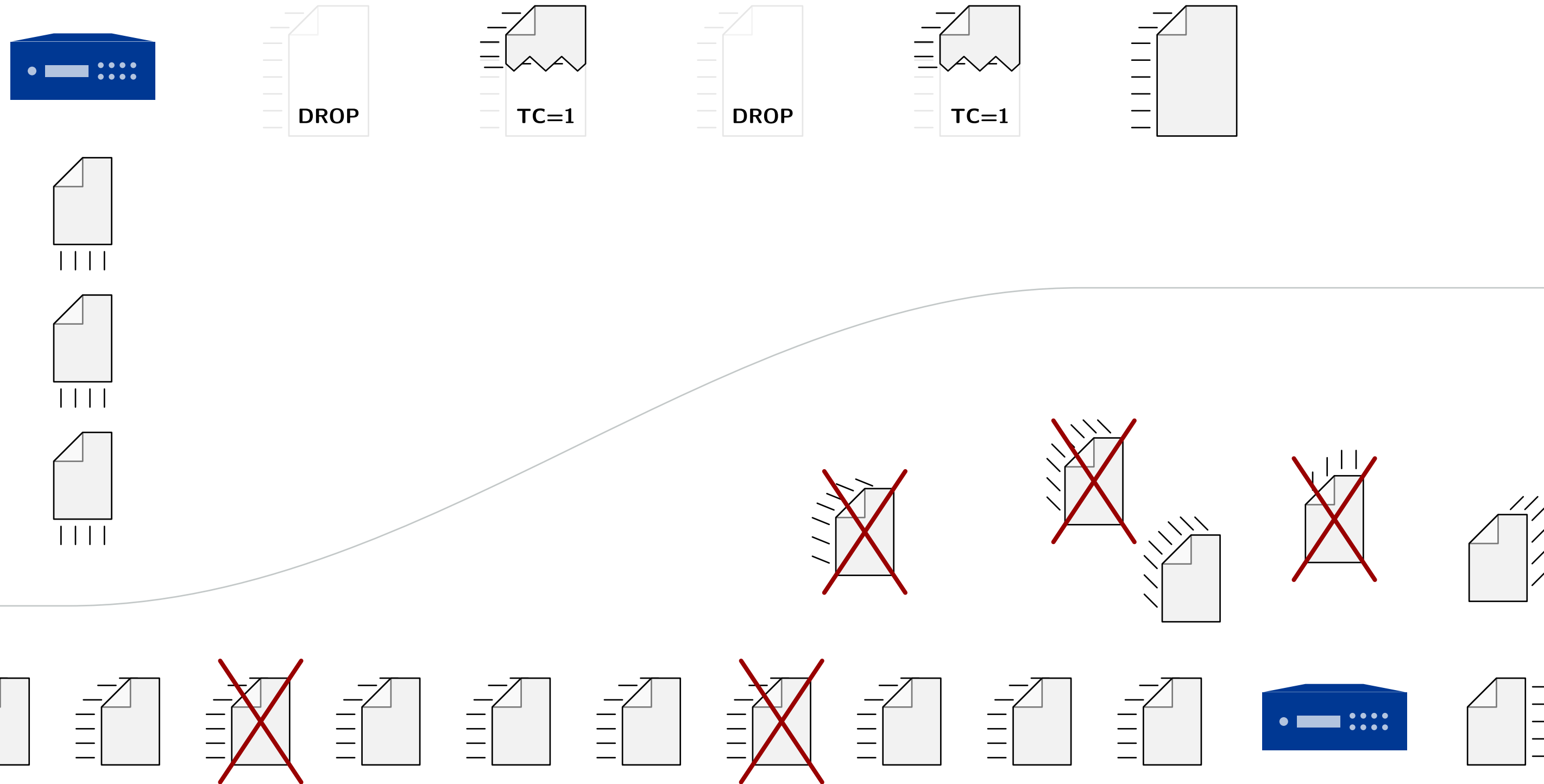


2. CPU overload

Overview

1. amplification att.

- forged UDP source
- answers >> queries
- size limit (1232 B)
- rate limiting
 - truncation
 - dropping



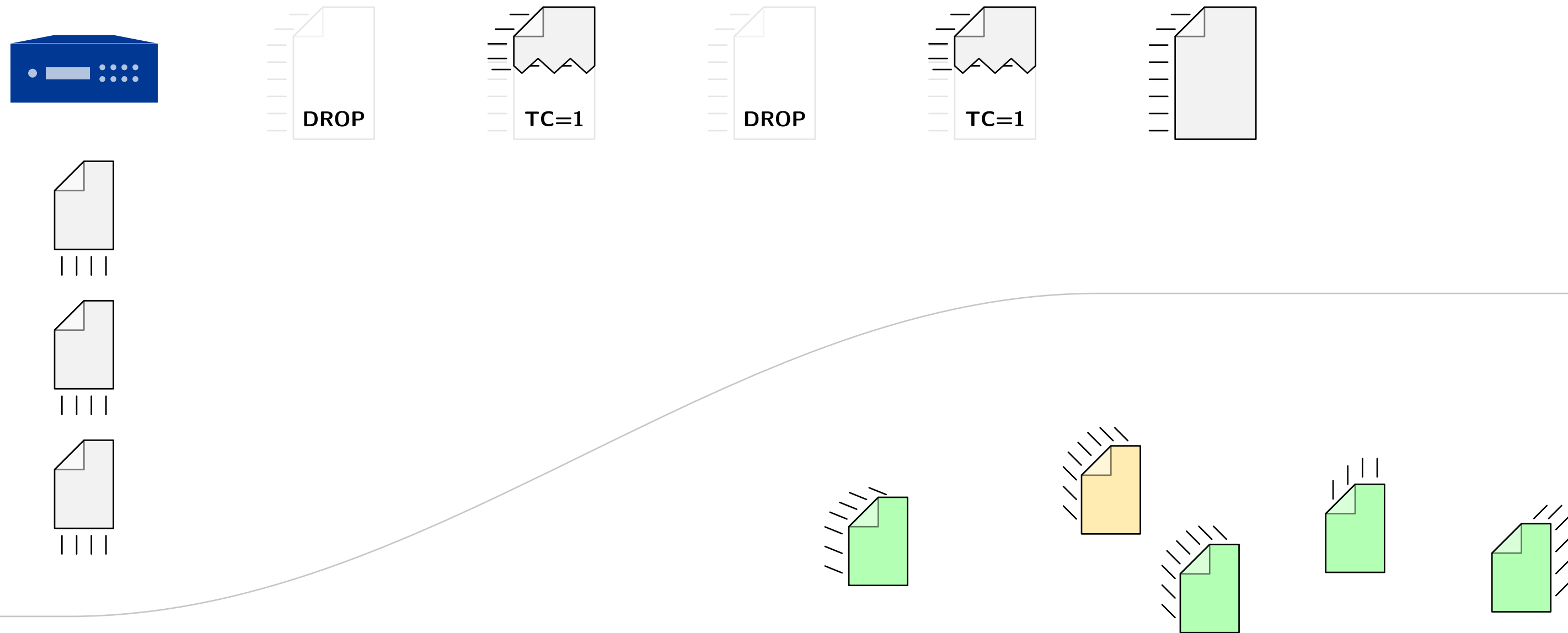
2. CPU overload

Overview

- Solution: defer requests from the origins using more cpu time in the past, so that users that do not overload the service shouldn't suffer.
- Non-UDP only, because on UDP the source IP could be faked.

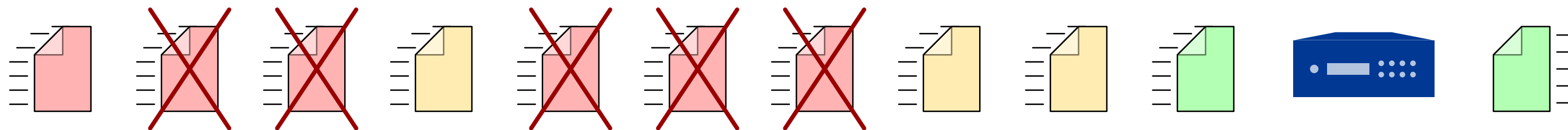
1. amplification att.

- forged UDP source
- answers >> queries
- size limit (1232 B)
- rate limiting
 - truncation
 - dropping



2. CPU overload

- prioritization
 - all except plain UDP

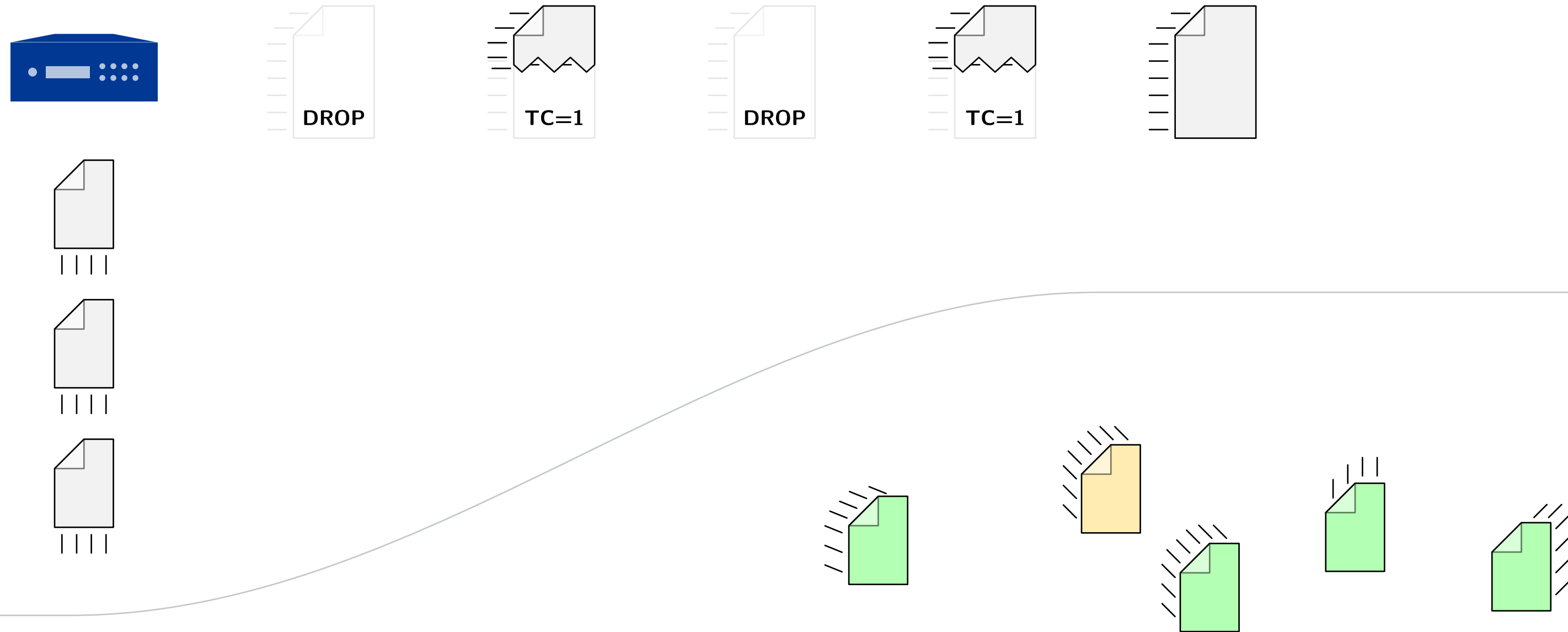


Overview

- Outline:
 - limiting individual hosts,
 - extending to networks,
 - different limits for dropping and truncating,
 - prioritization,
 - implementation details.

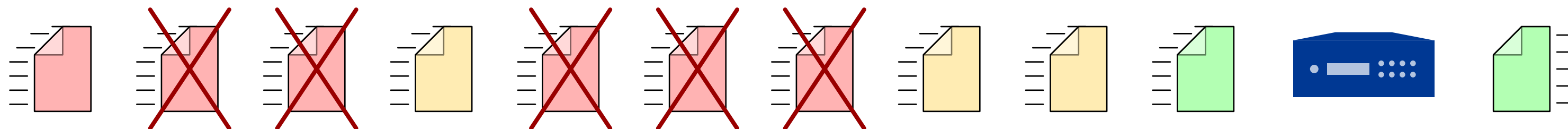
1. amplification att.

- forged UDP source
- answers >> queries
- size limit (1232 B)
- rate limiting
 - truncation
 - dropping



2. CPU overload

- prioritization
 - all except plain UDP



Limiting individual clients

- Mapping of addresses to counters – simplified.
- Values are counts of unrestricted queries, up to *instant* limit – max number of queries in a short period of time.

- **counters for addresses**

- **instant limit L_I**

	:
172.16.96.1	count in $[0, L_I)$
	⋮
2001:db8::734	count in $[0, L_I)$
	⋮

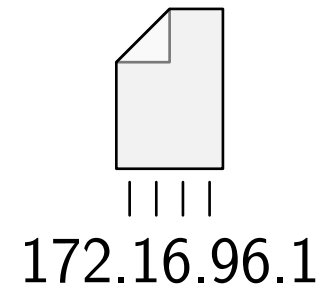
Limiting individual clients

- counters for addresses

- instant limit L_I

⋮	
172.16.96.1	count in $[0, L_I)$
⋮	
2001:db8::734	count in $[0, L_I)$
⋮	

← +1 unless overflows

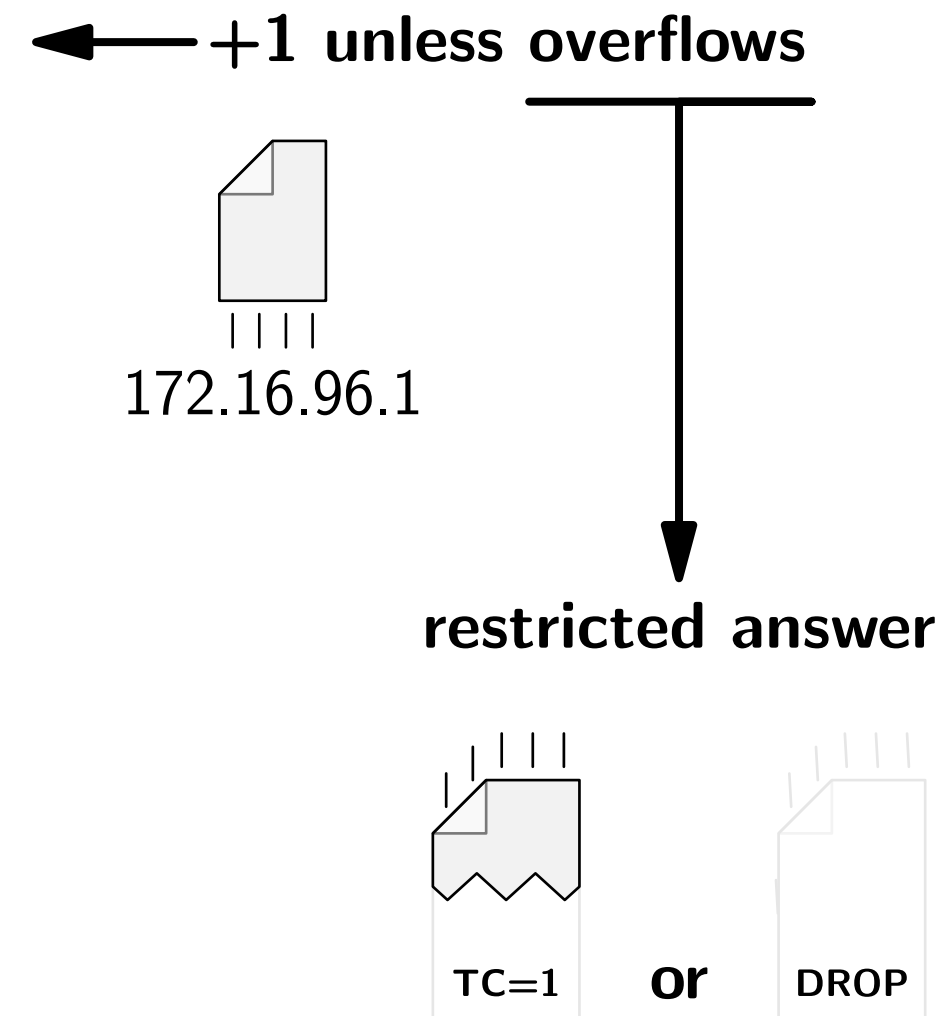


Limiting individual clients

- counters for addresses

- instant limit L_I

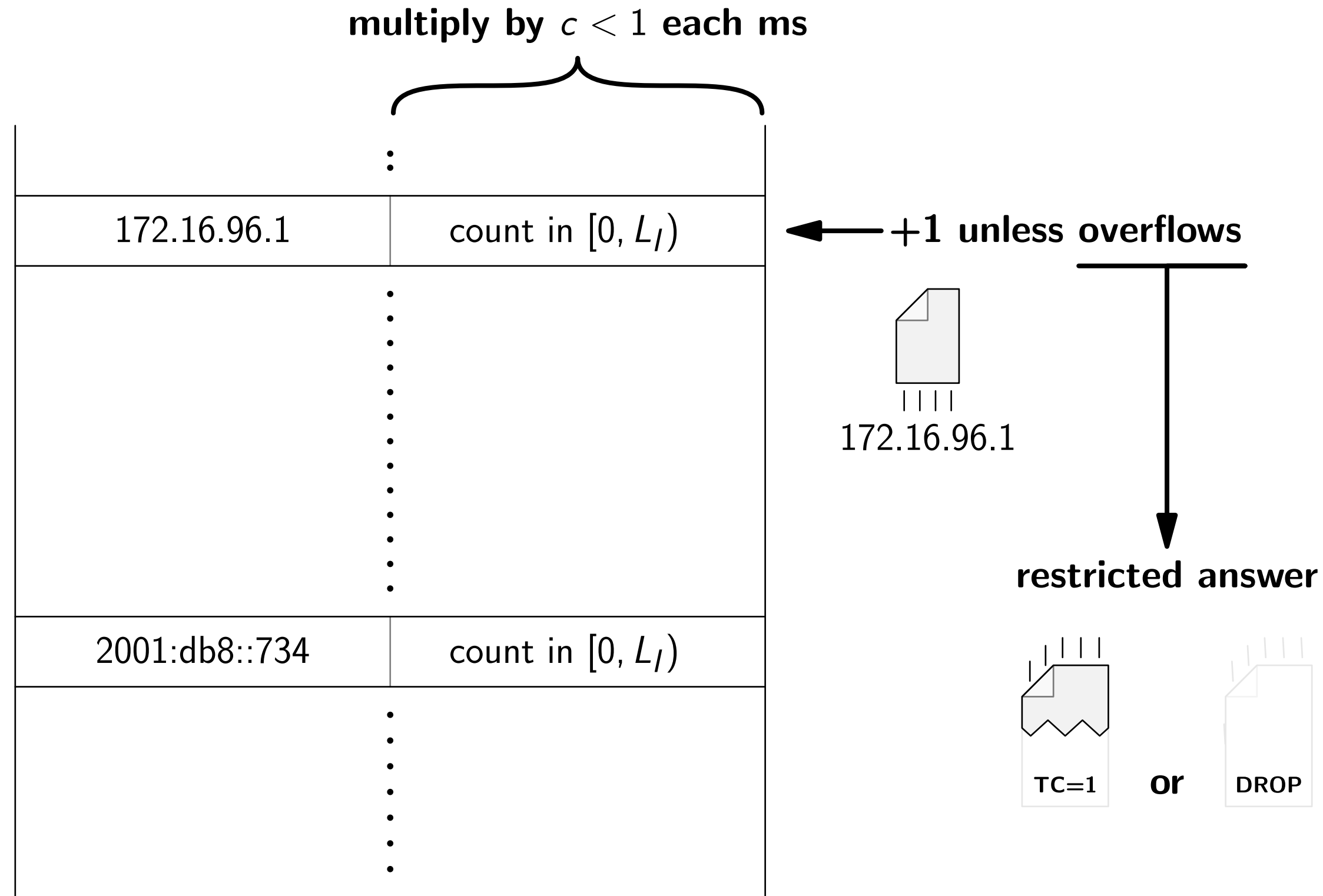
⋮	
172.16.96.1	count in $[0, L_I)$
⋮	
2001:db8::734	count in $[0, L_I)$
⋮	



Limiting individual clients

- Decreasing by a constant fraction of its value each ms.
- Exponential decay, resembles radioactive decay.
- The speed of decreasing given by *rate limit* – allowed number of queries per unit of time in the long-term.

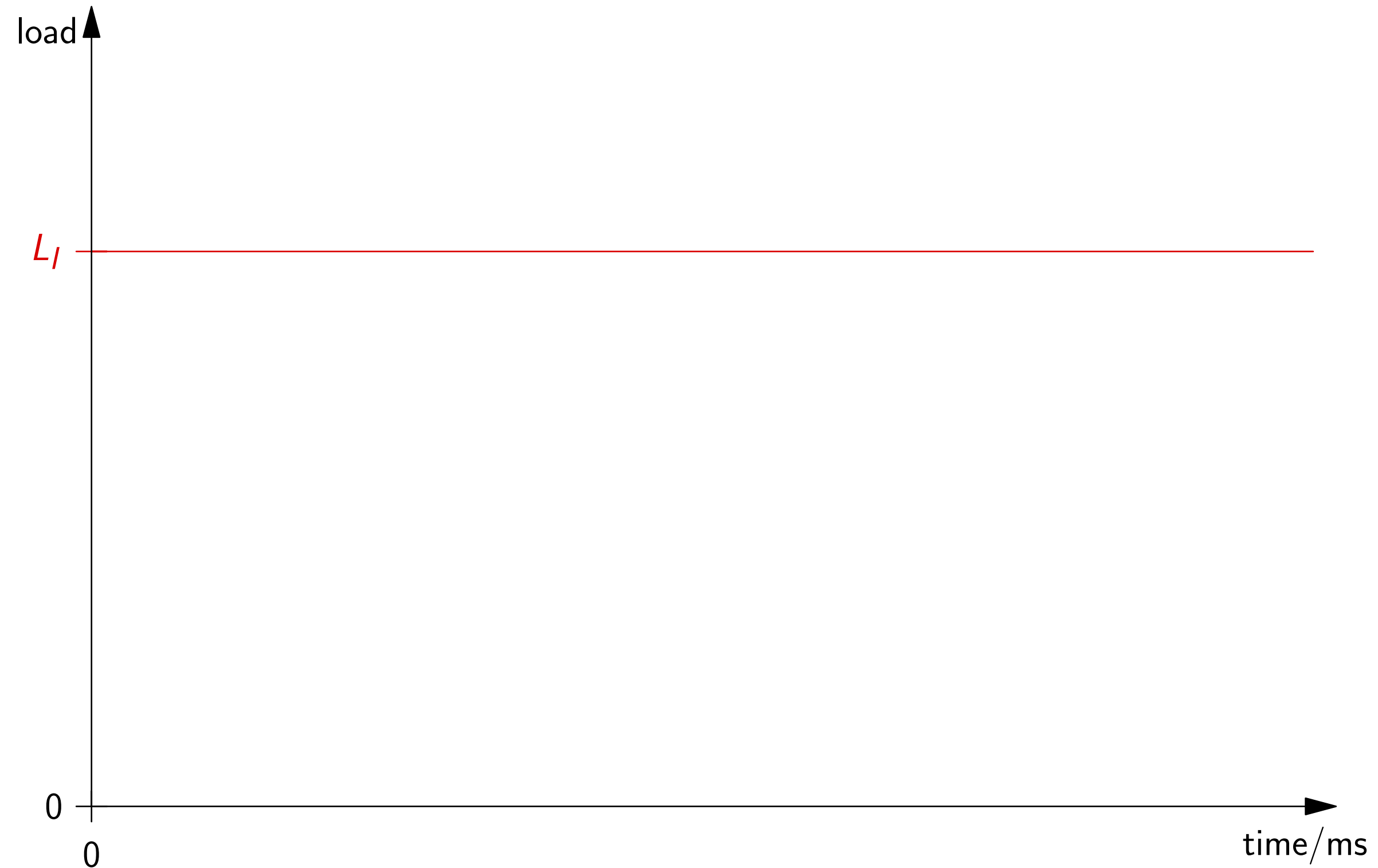
- counters for addresses
 - instant limit L_I
- exponential decay
 - rate limit



Exponential decay

- The limiting is configured by two values: Instant limit and (long-term) Rate limit.
- instant limit – max value of the counter

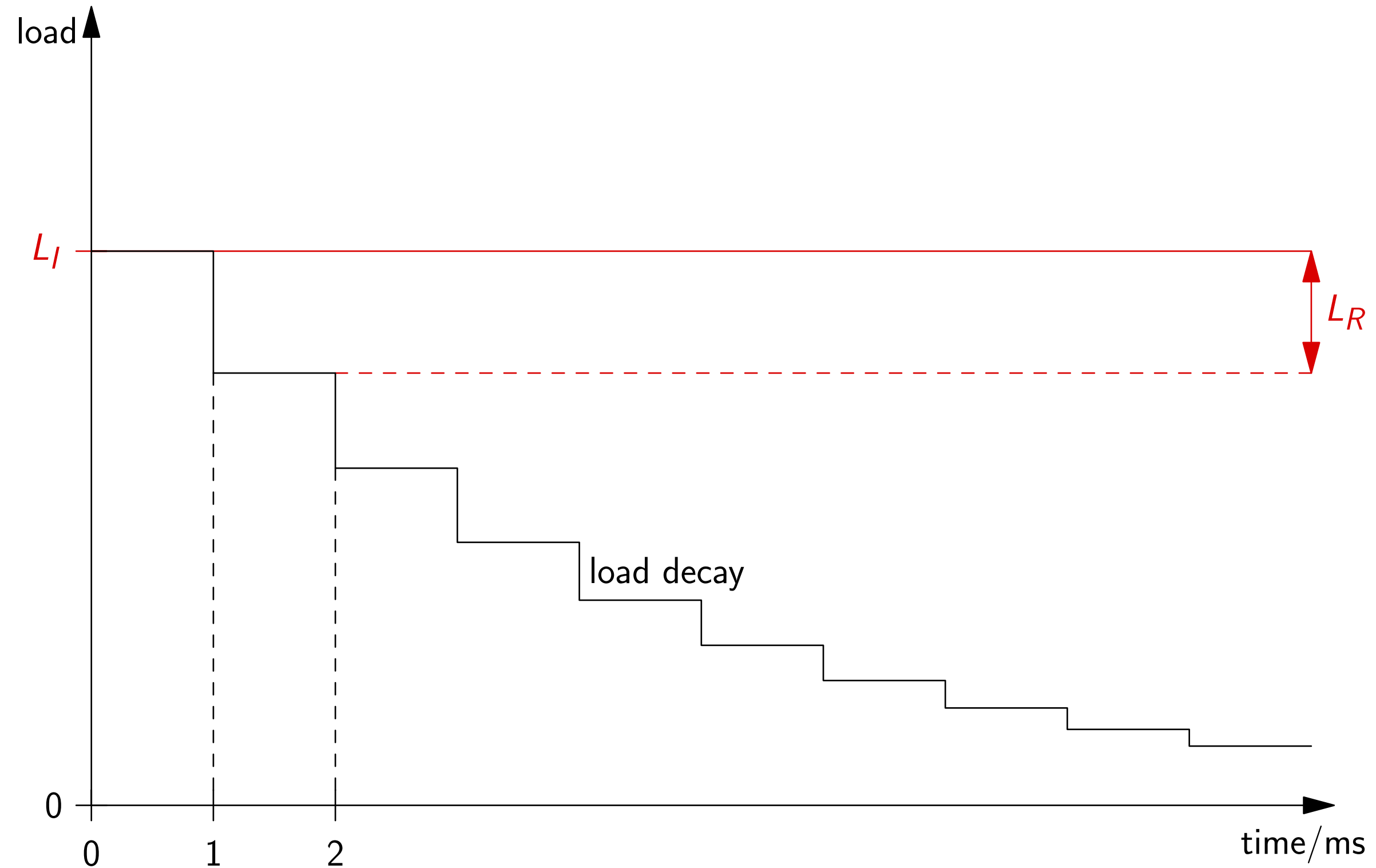
- instant limit L_I



Exponential decay

- decay after filling the counter
 - decreasing by constant fract. of the value each ms
- rate in ms – size of the first step
 - other steps lower

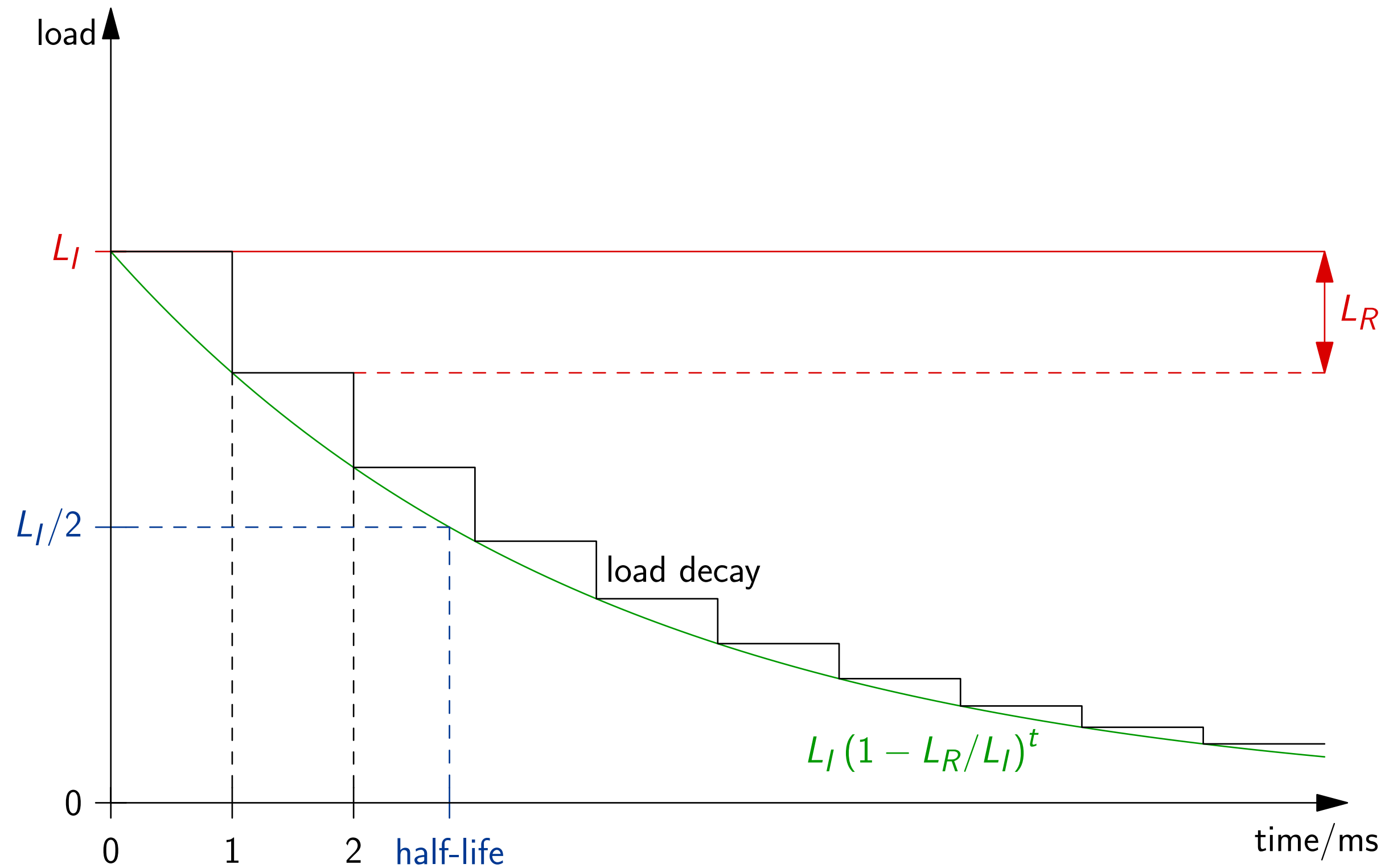
- instant limit L_I
- rate limit L_R
 - per ms



- (1/3)
- ratio of rate to instant gives half-life

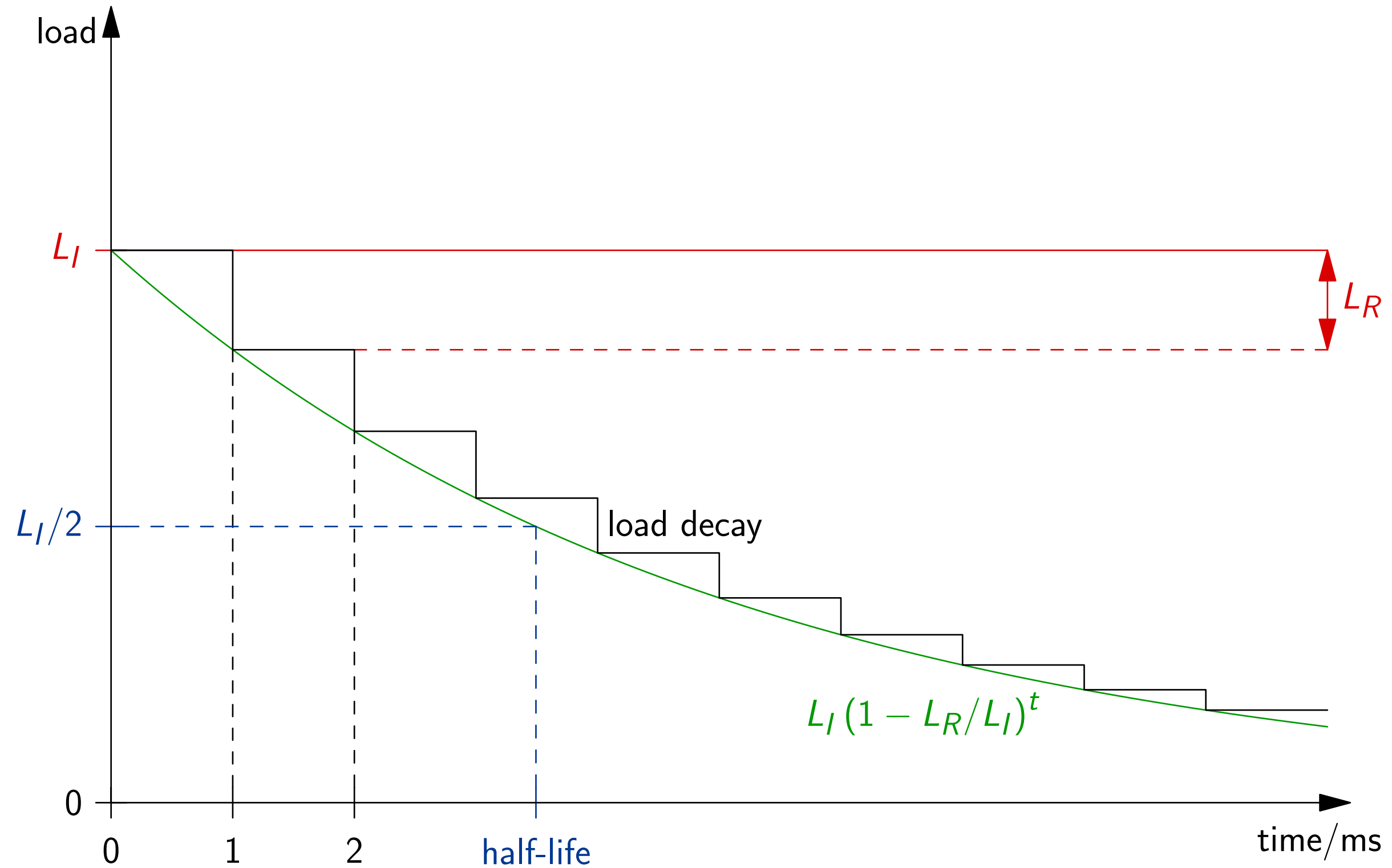
Exponential decay

- instant limit L_I
- rate limit L_R
 - per ms
- **half-life**



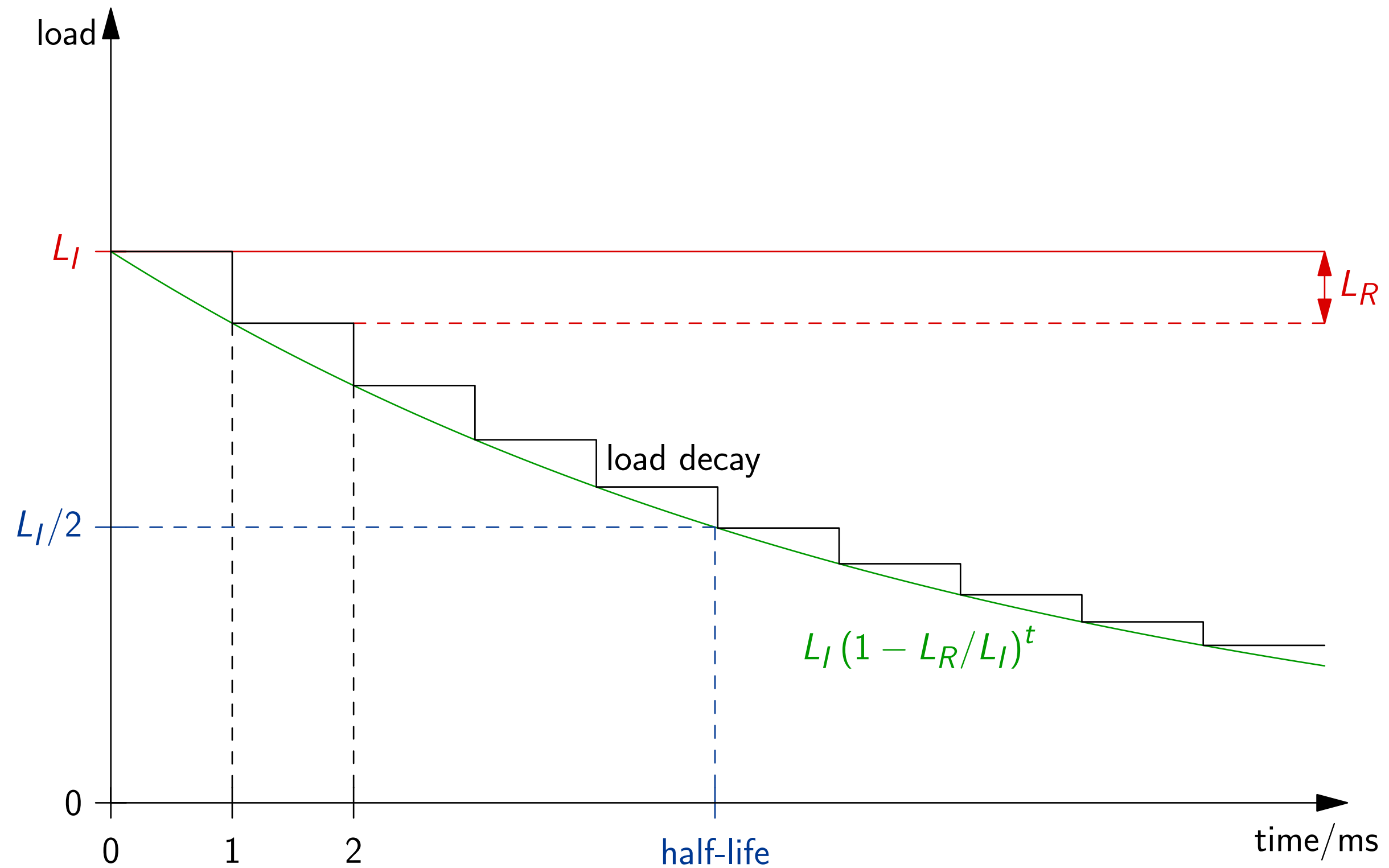
Exponential decay

- instant limit L_I
- rate limit L_R
 - per ms
- half-life



Exponential decay

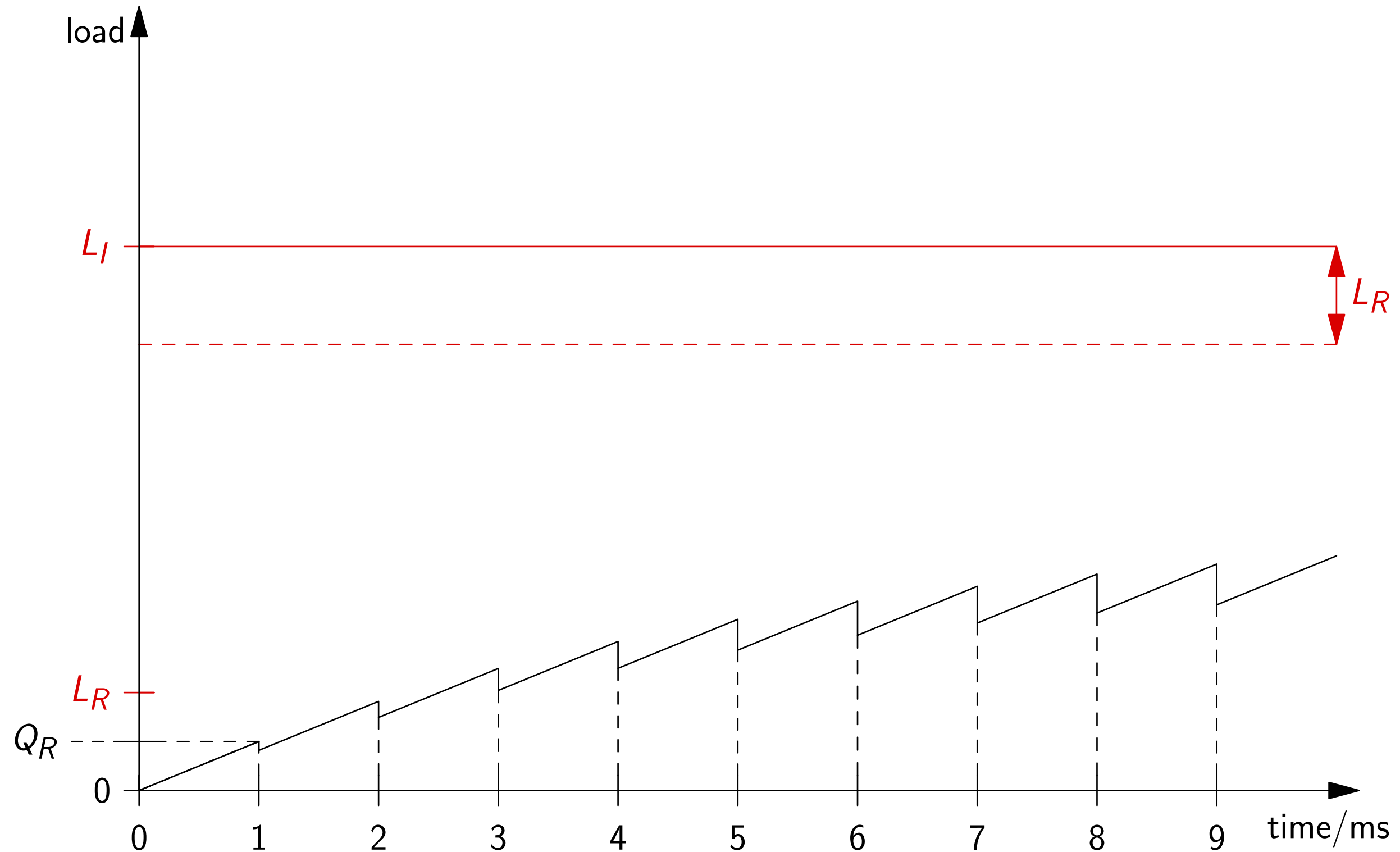
- instant limit L_I
- rate limit L_R
 - per ms
- half-life



Constant query rate example

- (1/3)
- Example: constant query rate under rate limit
 - no restriction

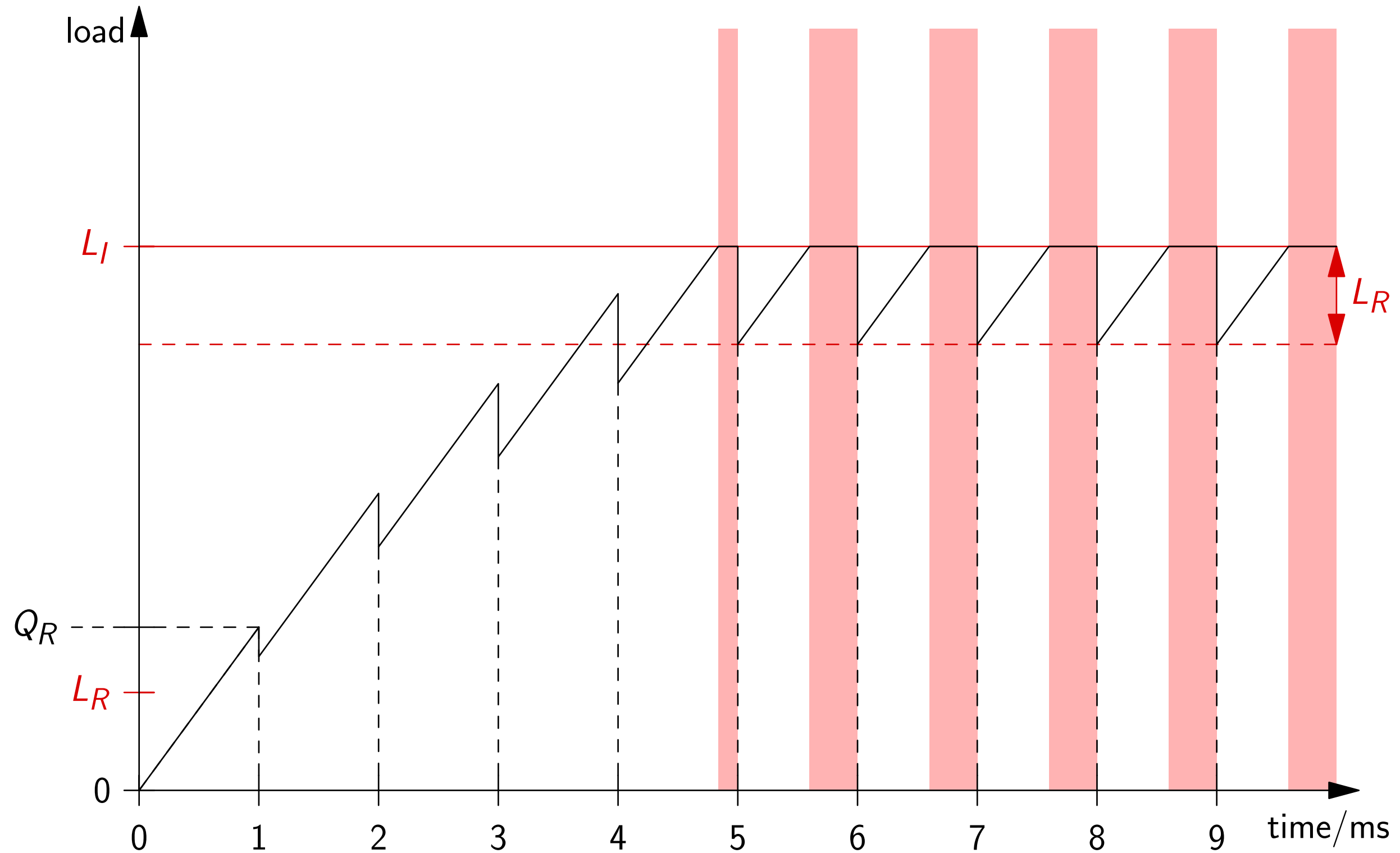
- instant limit L_I
- rate limit L_R
 - per ms
- half-life
- query rate Q_R
 - per ms



Constant query rate example

- (2/3)
- query rate above rate limit
 - gets restricted at some point in time

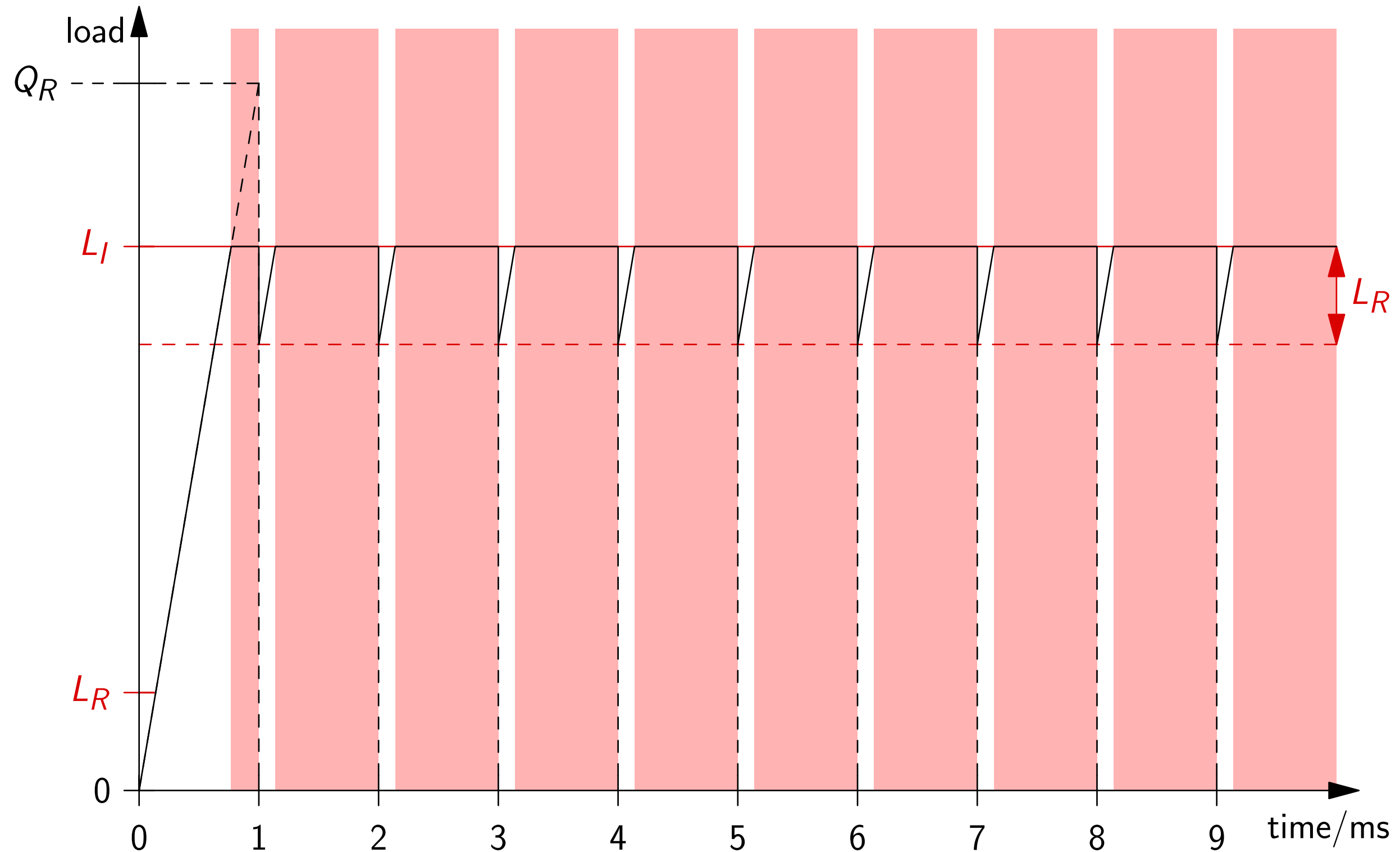
- instant limit L_I
- rate limit L_R
 - per ms
- half-life
- query rate Q_R
 - per ms



Constant query rate example

- (3/3)
- query rate above instant limit
 - restricted just after reaching L_I
 - then according to rate limit
- unrestricted queries -> something passes through, response rate is limited

- instant limit L_I
- rate limit L_R
 - per ms
- half-life
- query rate Q_R
 - per ms



Limiting networks

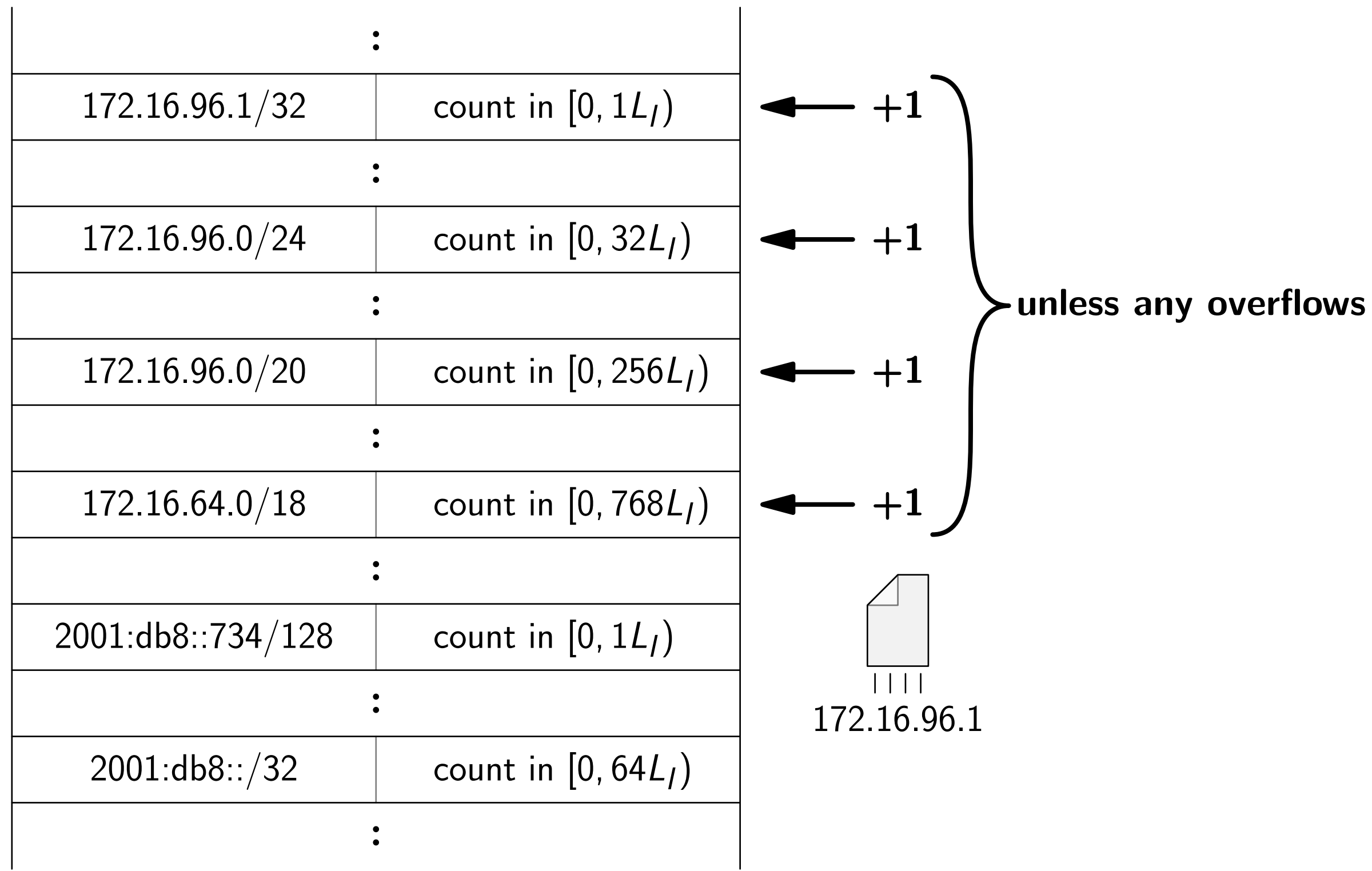
- Attackers won't play nice.
- Limiting individual IPs doesn't suffice.
- Usual approach: limit a single prefix size.
- We: maintain counters for *several* chosen prefixes.
- Constants multiplying limits based on prefix size.
- Shorter prefix -> larger network -> higher limits.
- Same multiplier for rate and instant -> half-life unchanged.

- IPv4
 - /32: 1
 - /24: 32
 - /20: 256
 - /18: 768
- IPv6
 - /128: 1
 - /64: 2
 - /56: 3
 - /48: 4
 - /32: 64

	:
172.16.96.1/32	count in $[0, 1L_l)$
	:
172.16.96.0/24	count in $[0, 32L_l)$
	:
172.16.96.0/20	count in $[0, 256L_l)$
	:
172.16.64.0/18	count in $[0, 768L_l)$
	:
2001:db8::734/128	count in $[0, 1L_l)$
	:
2001:db8::/32	count in $[0, 64L_l)$
	:

Limiting networks

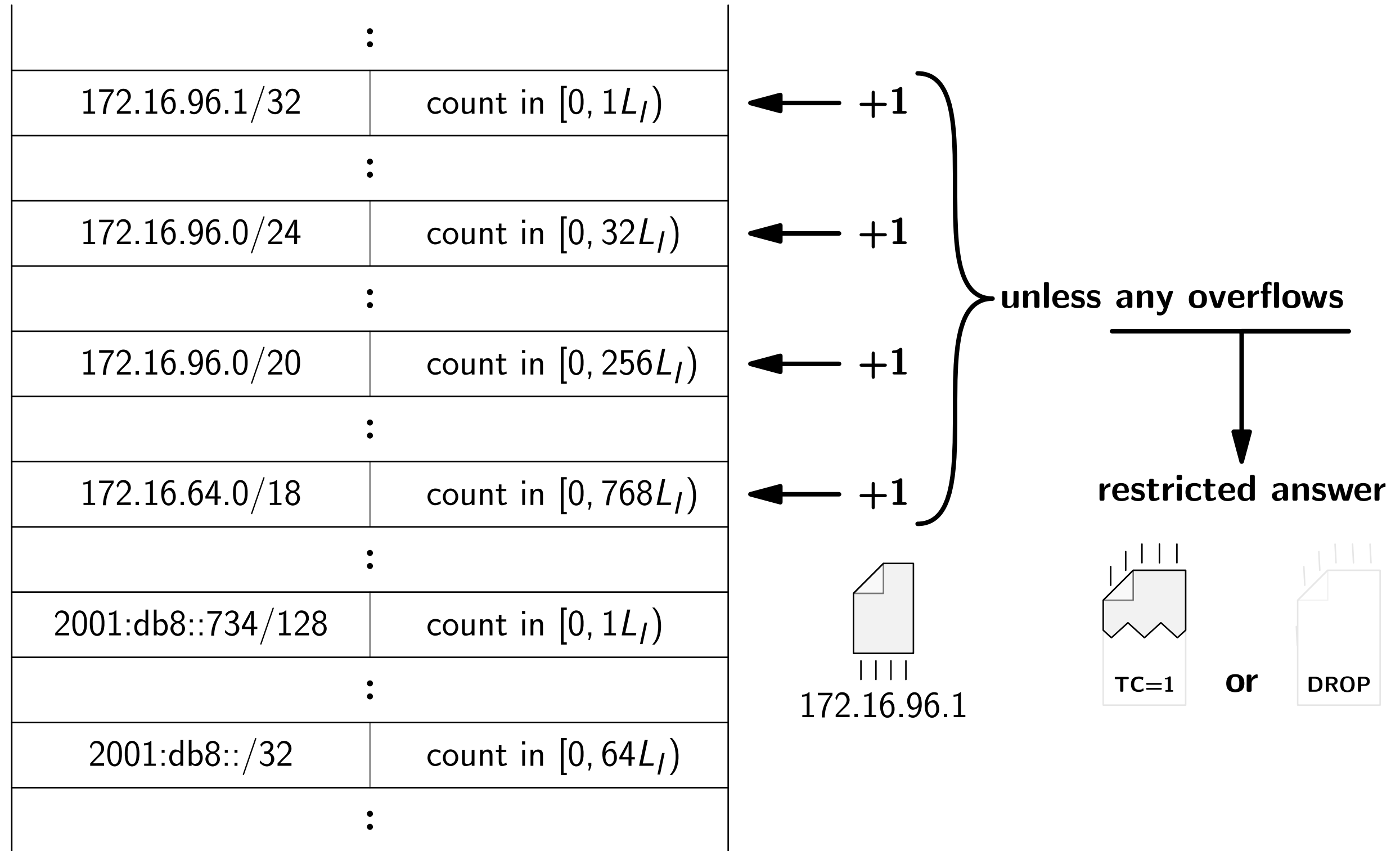
- IPv4
 - /32: 1
 - /24: 32
 - /20: 256
 - /18: 768
- IPv6
 - /128: 1
 - /64: 2
 - /56: 3
 - /48: 4
 - /32: 64



- ...in which case restricted.
 - read-only, faster

Limiting networks

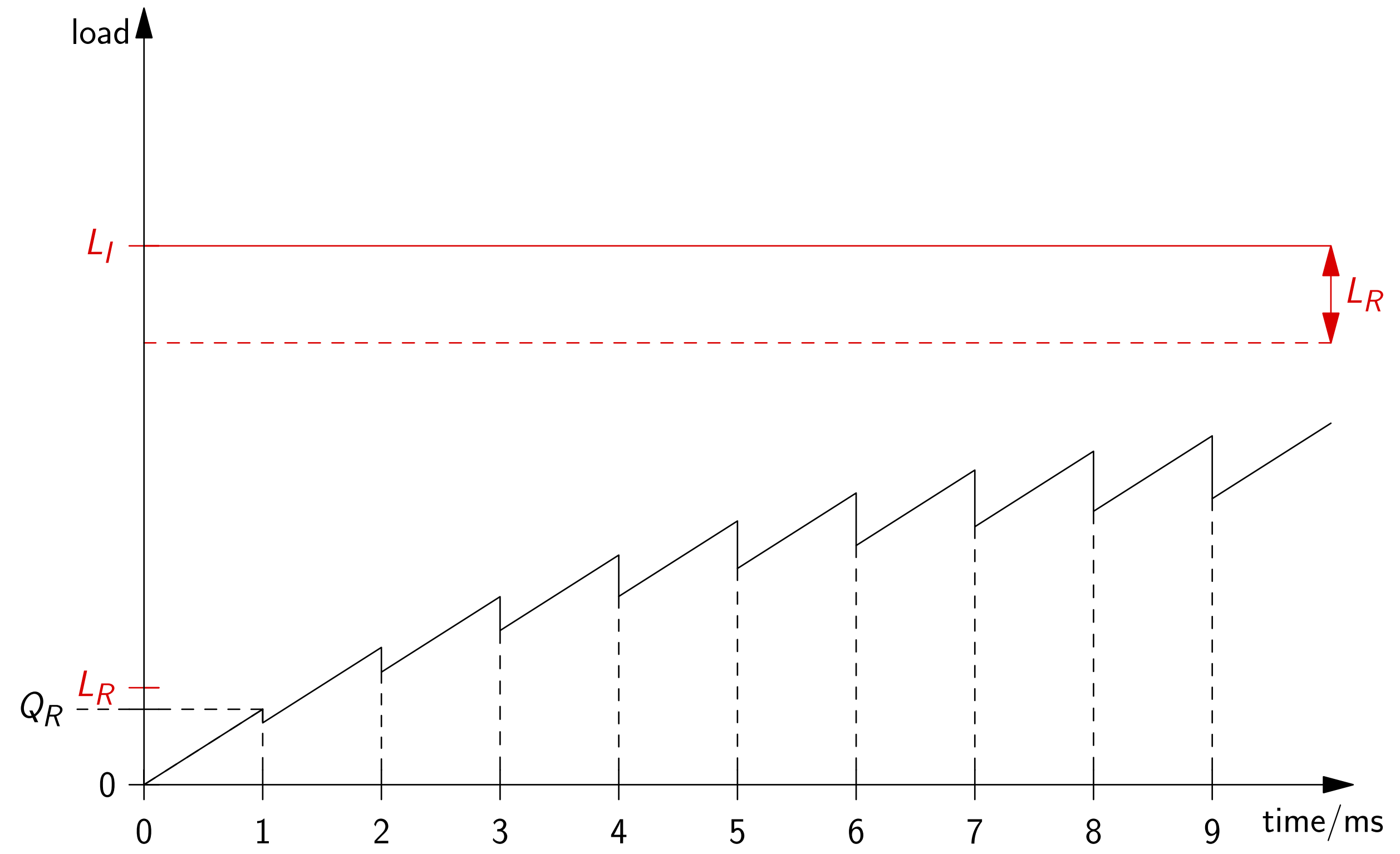
- IPv4
 - /32: 1
 - /24: 32
 - /20: 256
 - /18: 768
- IPv6
 - /128: 1
 - /64: 2
 - /56: 3
 - /48: 4
 - /32: 64



- So far only hard limits for dropping.
- Add lower instant and rate limits for truncating.

Multiple limits

- hard limits
 - instant L_I
 - rate L_R per ms



Multiple limits

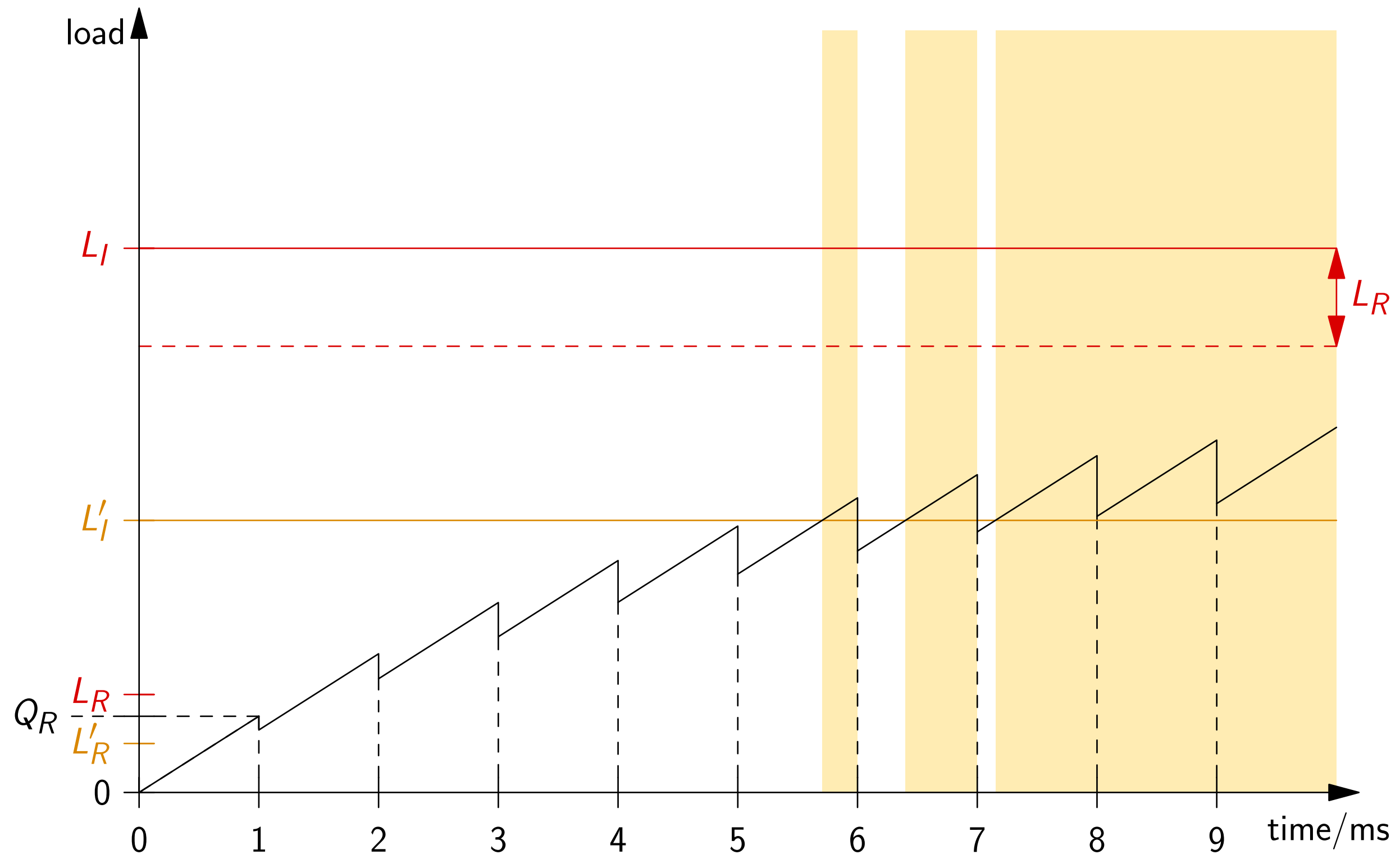
- (1/2)
- Incrementing also over soft limit – otherwise cannot reach hard limit.
- Everything truncated until user lowers query rate.
- On the chart query rate between soft and hard rate limit.

- **hard limits**

- instant L_I
- rate L_R per ms

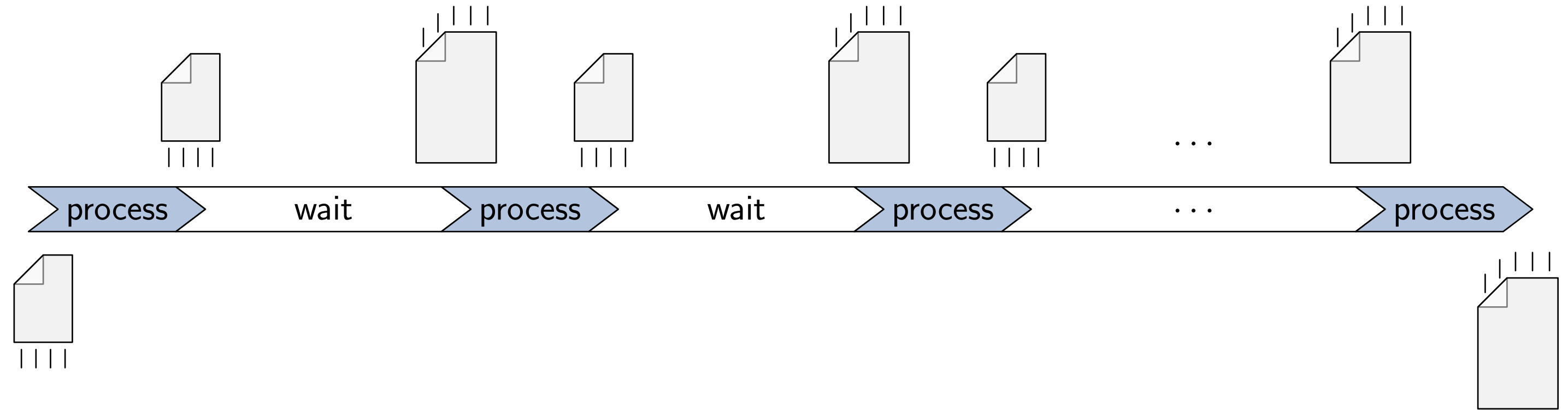
- **soft limits**

- instant L'_I
- rate L'_R per ms



Prioritization

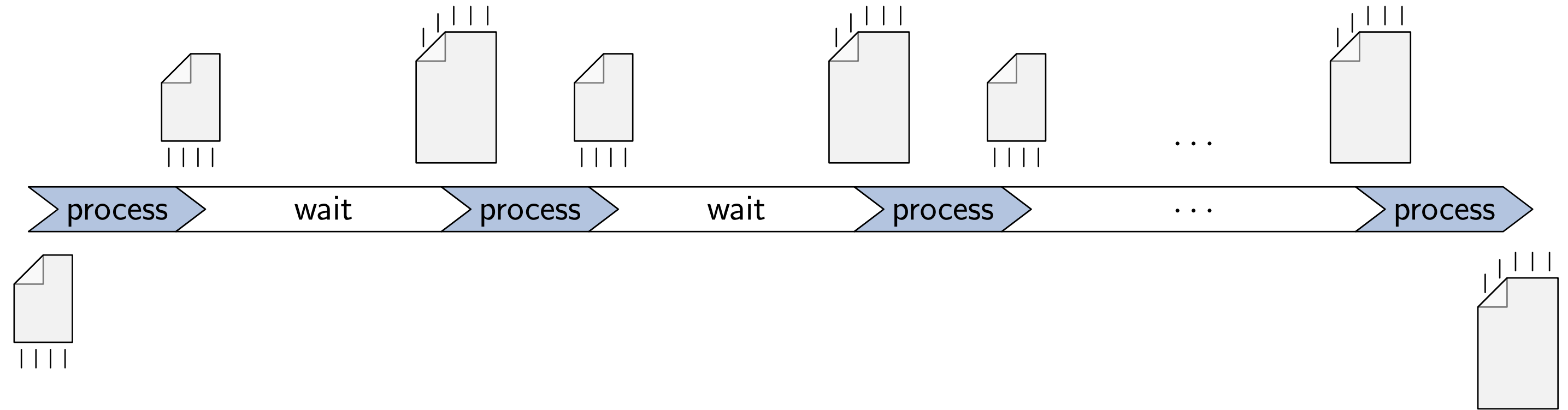
- not on plain UDP
- no configuration



Prioritization

- Even in legitimate traffic some queries are way more expensive than average.
- Aim: catch as much as possible while prioritizing users that don't overload our CPU.
- CPU time is measured, waiting not.

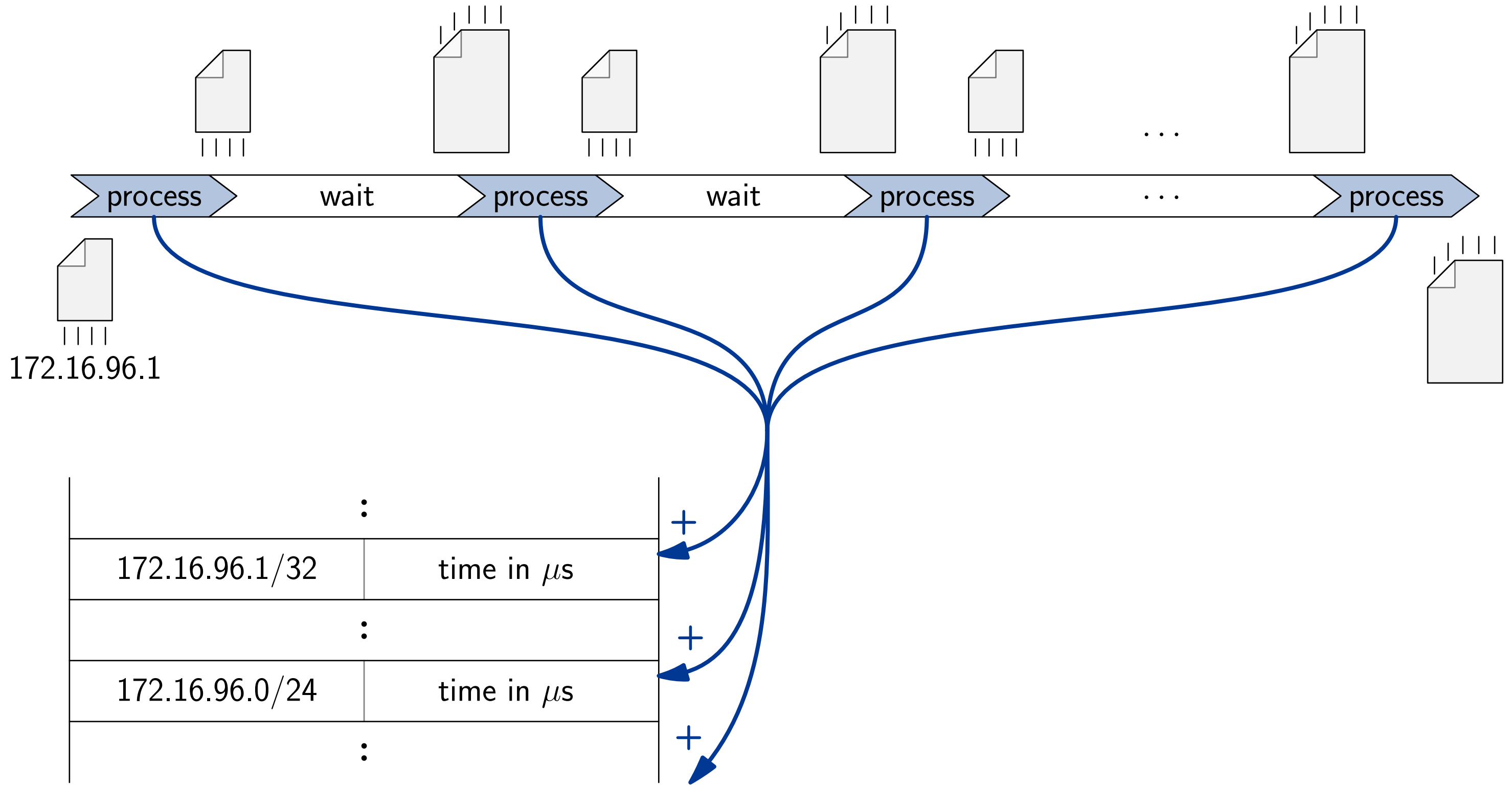
- not on plain UDP
- no configuration
- measuring time
 - only cpu, no wait



Prioritization

- Incrementing counters by time in μs .
 - different table instance
 - both addresses and networks

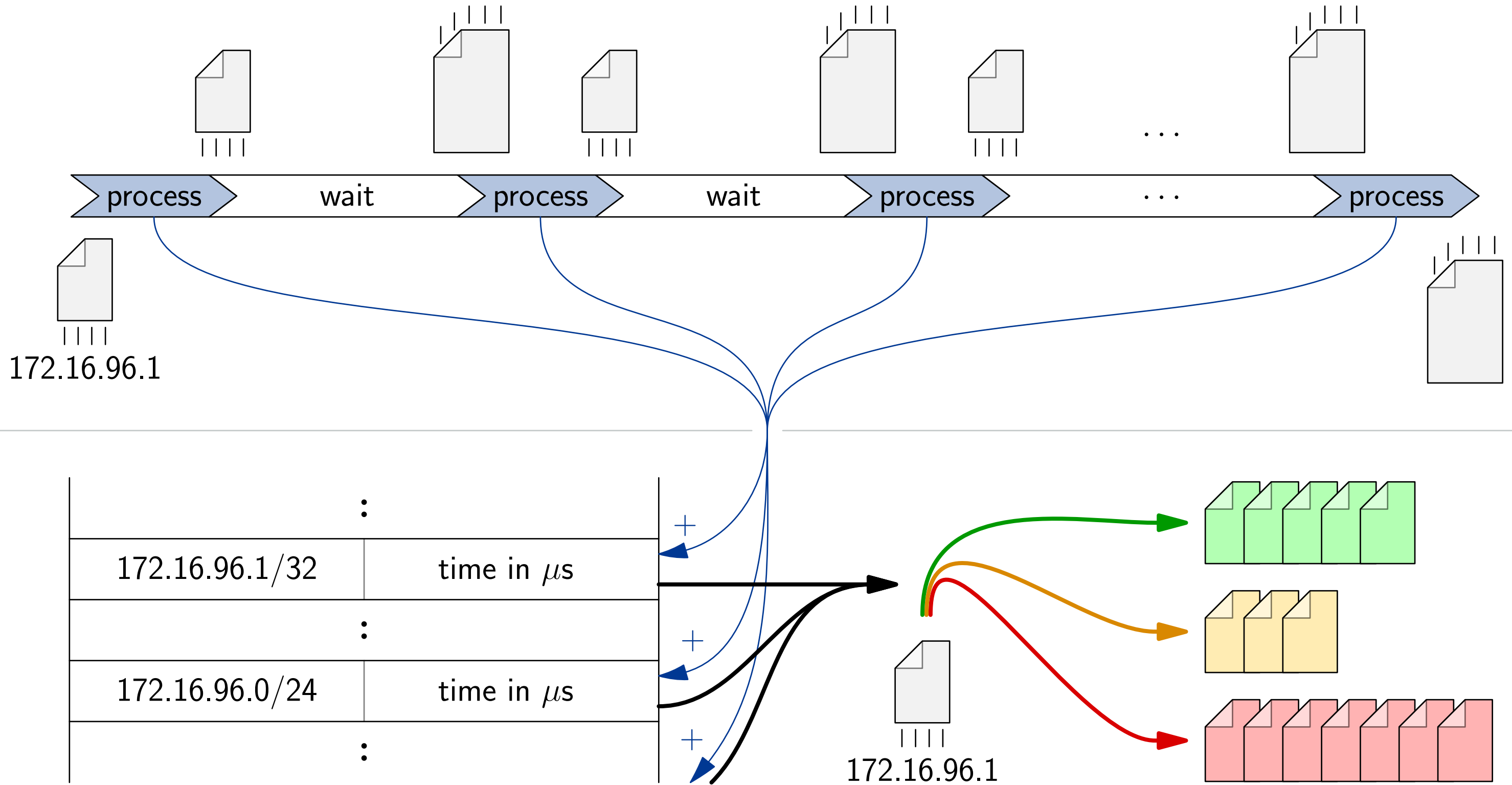
- not on plain UDP
- no configuration
- measuring time
 - only cpu, no wait
 - **add to table values**



Prioritization

- Multiple soft limits for different priorities.
 - A queue for each priority.
 - May be deferred multiple times on priority decrease.

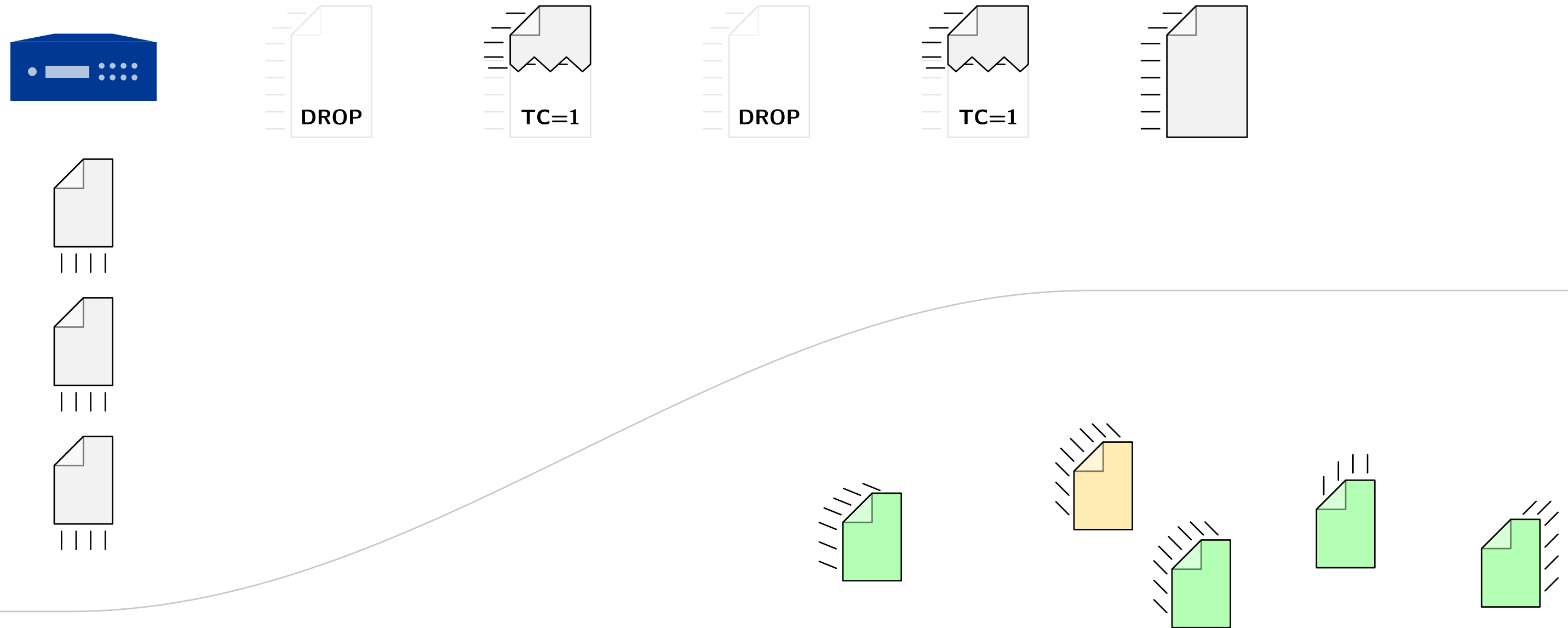
- not on plain UDP
- no configuration
- measuring time
 - only cpu, no wait
 - add to table values
- defer to queues
 - based on values



Final overview

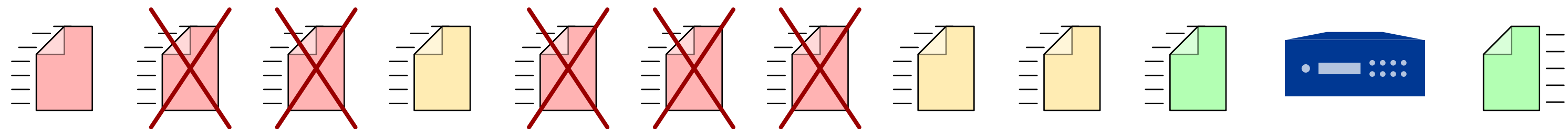
1. amplification att.

- forged UDP source
- answers >> queries
- size limit (1232 B)
- rate limiting
 - truncation
 - dropping



2. CPU overload

- prioritization
 - all except plain UDP



Appendix

now extra slides not planned for LinuxDays

Implementation

- Not possible to store all addresses, we store the most important ones (to be defined later).
- Use hash table to store them.
- Collisions may occur.

- **hashing**

hash(172.16.96.1/32) = 10111101100001101010110000110101001110100111010011101101001000001011111

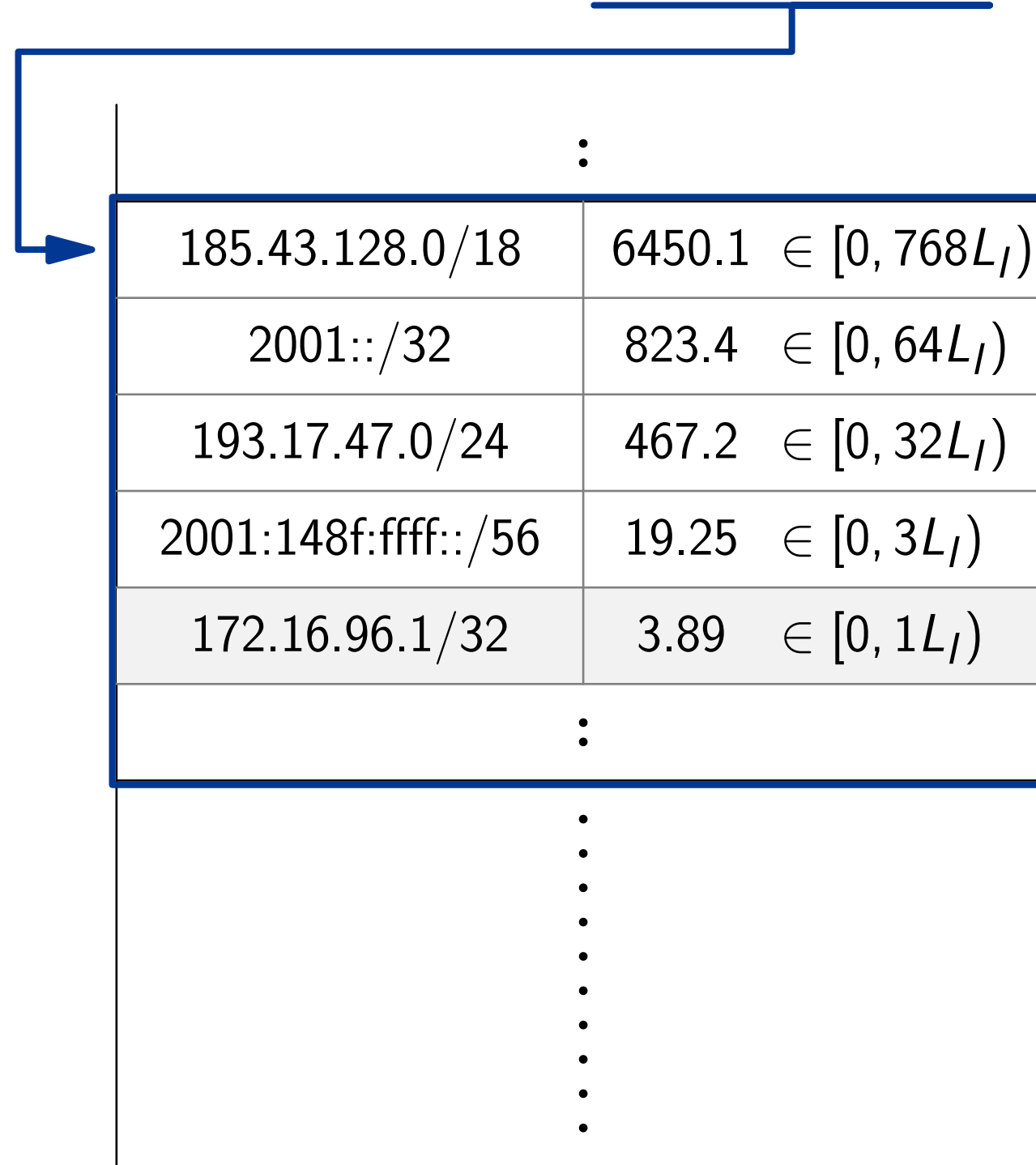
:	:
172.16.96.1/32	3 ∈ [0, 1L _I)
:	:
172.16.96.0/24	15.34 ∈ [0, 32L _I)
:	:
172.16.96.0/20	123 ∈ [0, 256L _I)
:	:
2001:db8::734/128	7.569 ∈ [0, 1L _I)
:	:
2001:db8::/32	33.21 ∈ [0, 64L _I)
:	:

Implementation

- Use buckets with several (15) most important records.
- Still low number of buckets will have many collisions.

- hashing
 - buckets

hash(172.16.96.1/32) = 10111101100001101010110000110101001110100111010011101101001000001011111

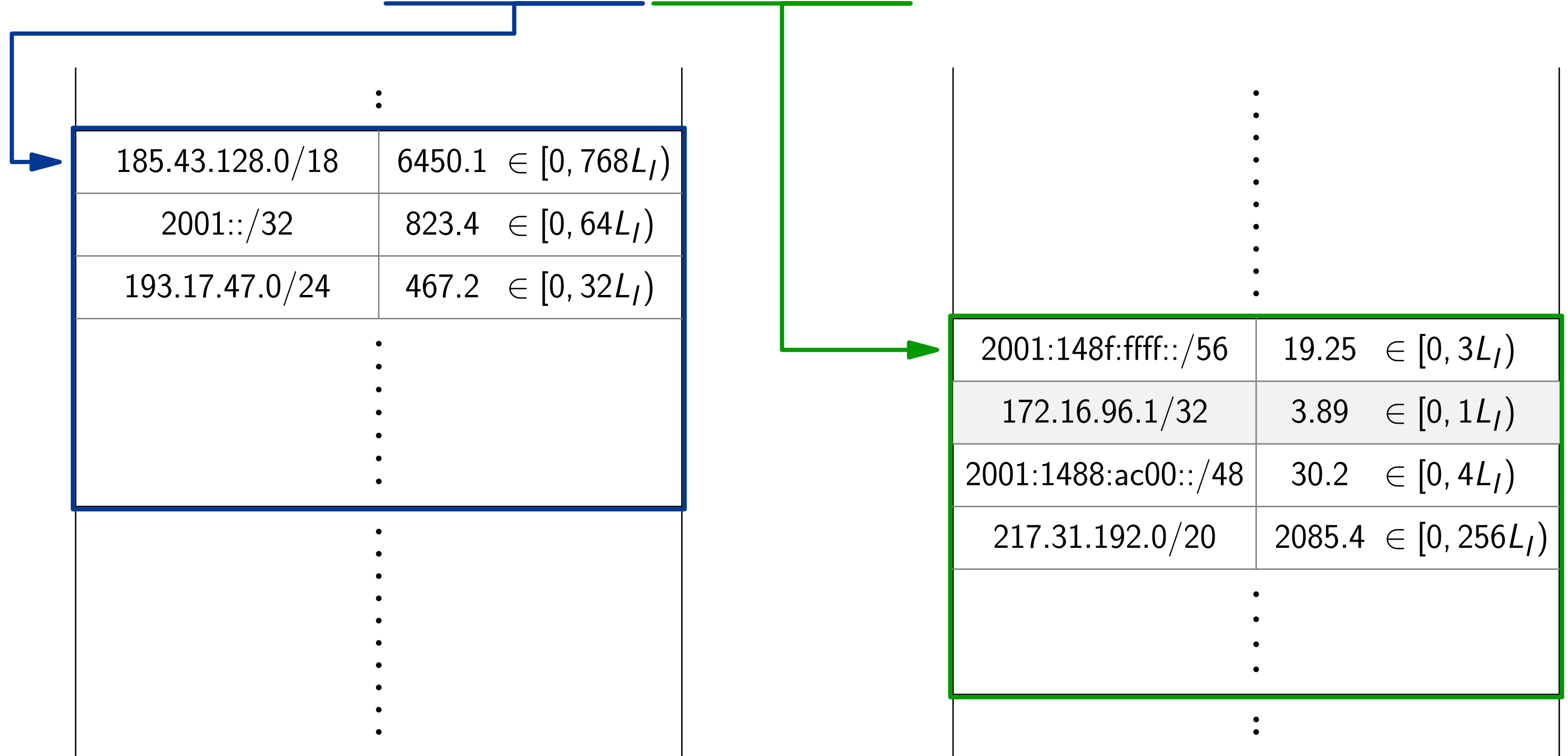


Implementation

- Two tables, hashed independently.
- The probability of collision in both of them is much smaller.

- hashing
 - buckets
 - **two tables**

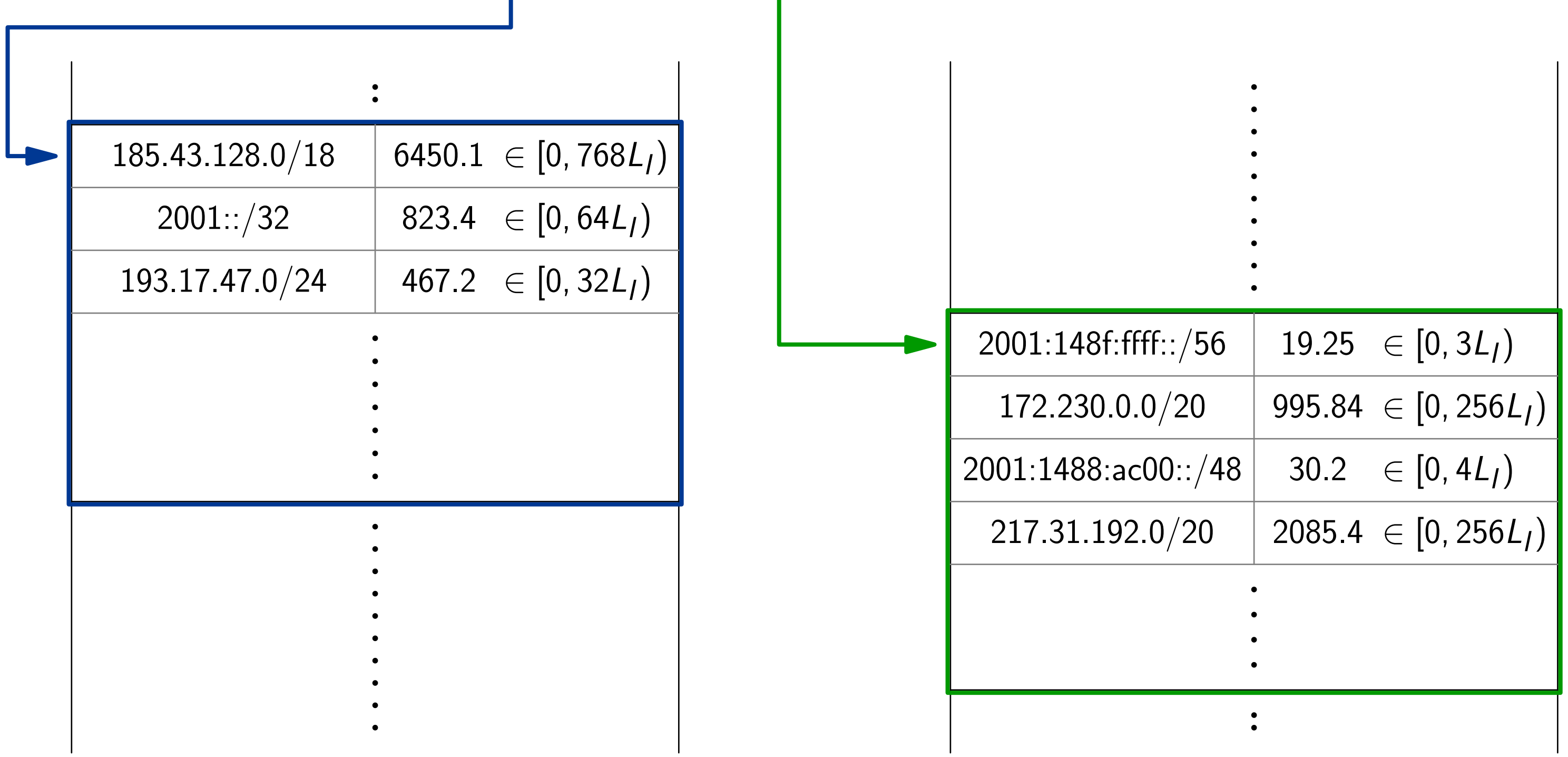
hash(172.16.96.1/32) = 1011110110000110101011000011010100111010011101101001000001011111



Implementation

- hashing
 - buckets
 - two tables
- evicting

hash(172.16.96.1/32) = 101111011000011010101100001101010011101001110100111011010010000010111111

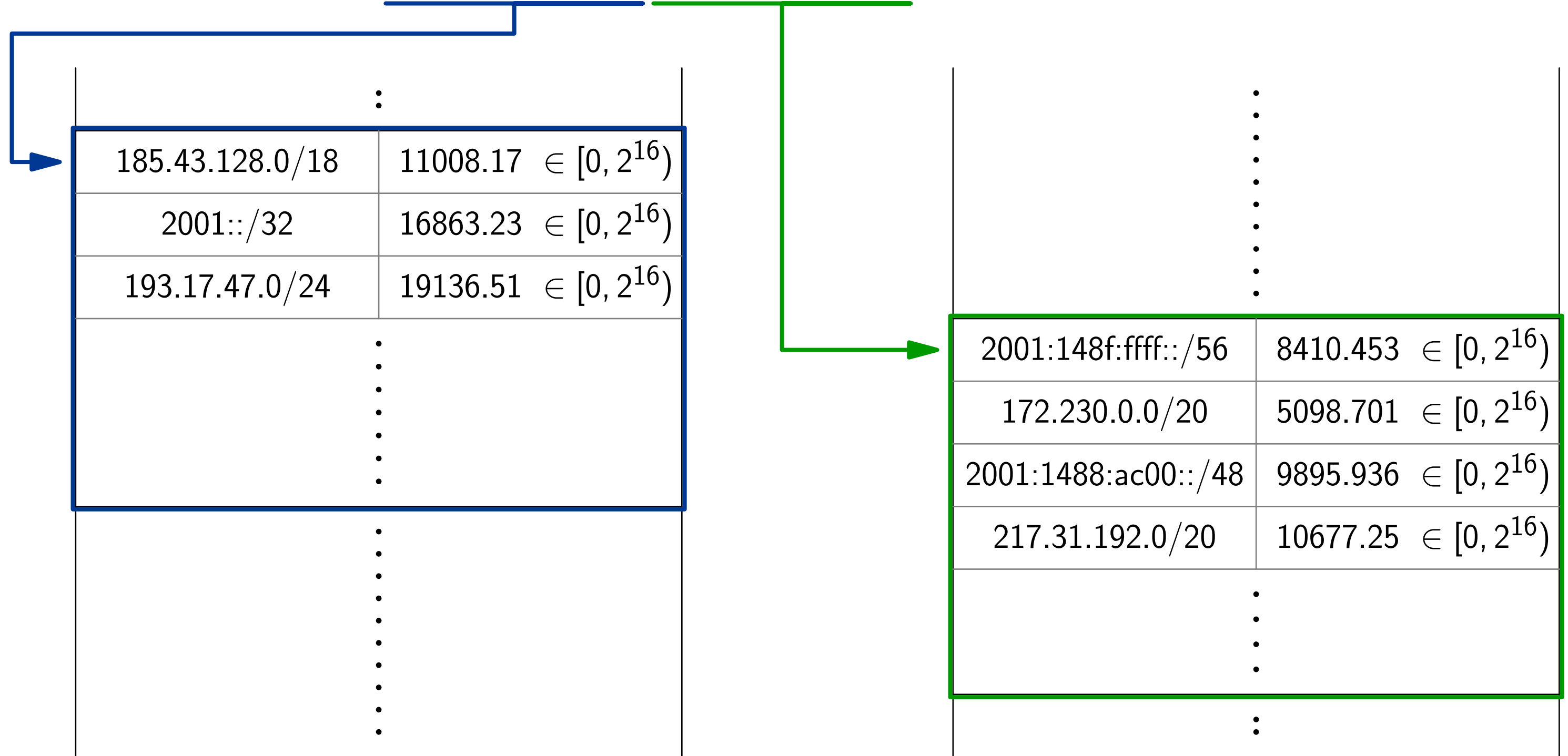


Implementation

- Normalize to the same limit.
- It allows comparing values across different prefix length – gives us the notion of their importance.

- hashing
 - buckets
 - two tables
- evicting
 - normalized limits

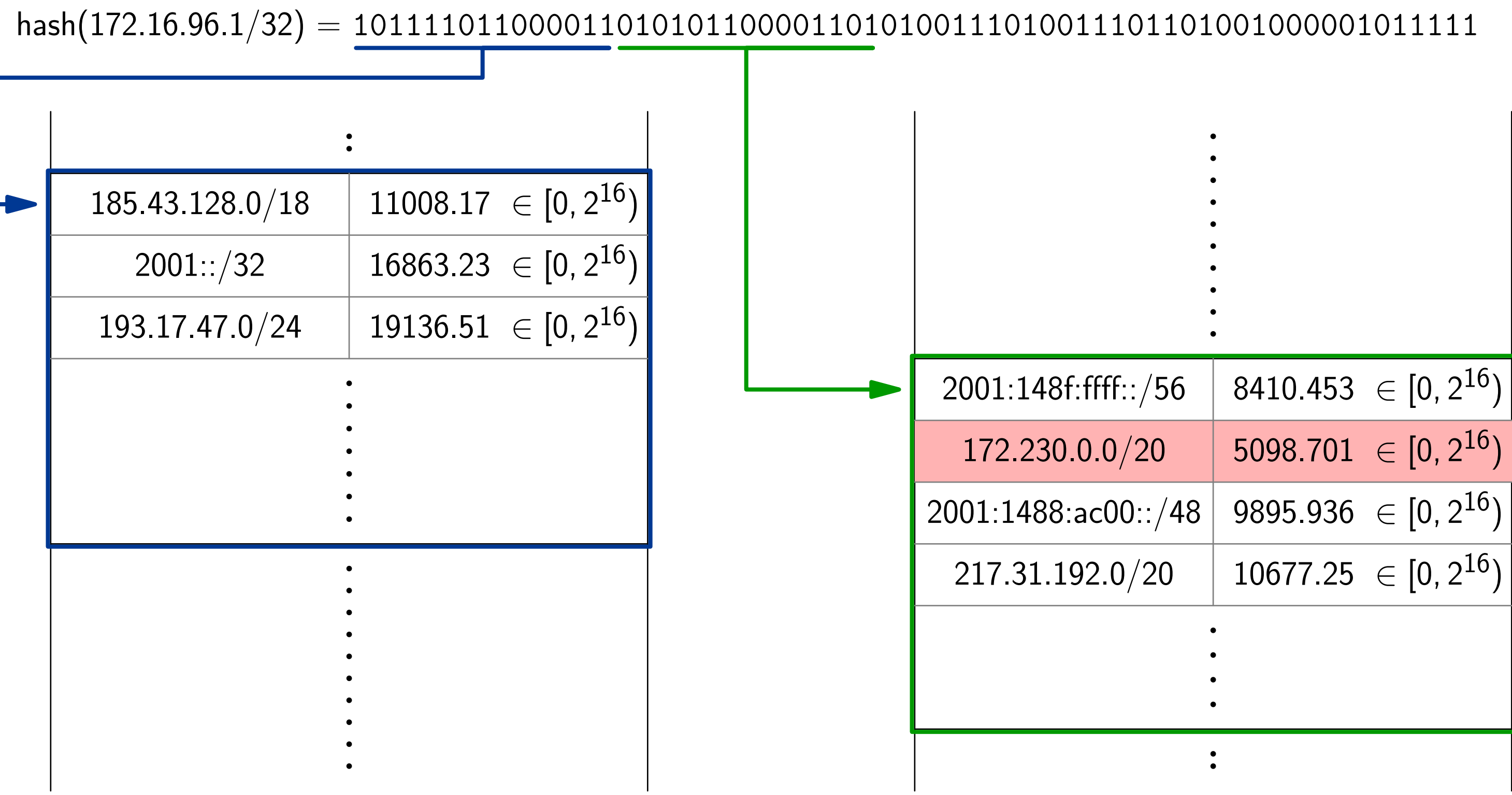
hash(172.16.96.1/32) = 10111101100001101010110000110101001110100111010011101101001000001011111



- If both buckets are full and new record appears, evict the one with the smallest value.

Implementation

- hashing
 - buckets
 - two tables
- evicting
 - normalized limits
 - **choosing minimal**

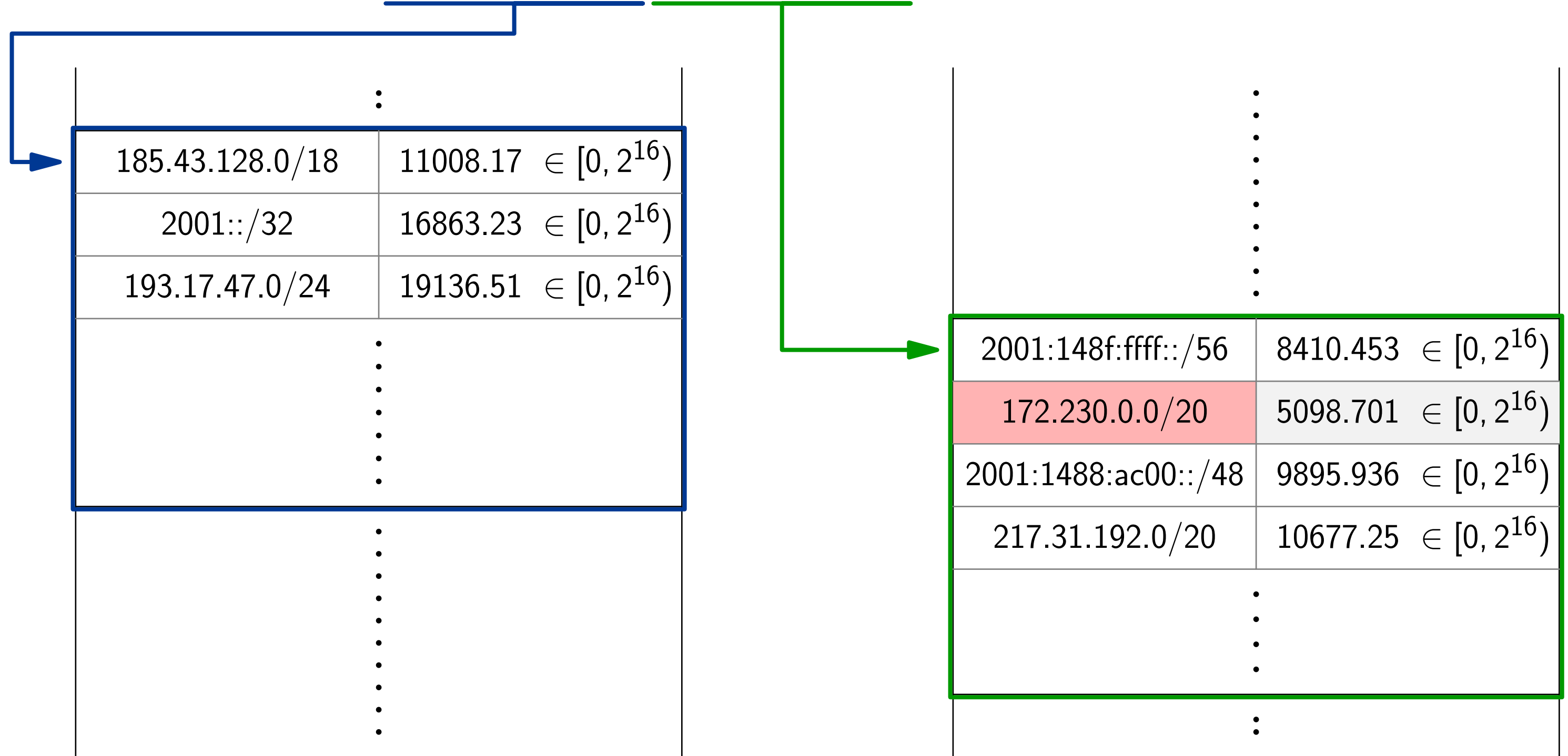


Implementation

- In fact, evict only the label keeping value.
- Multiple items evicting each other share the value instead of zeroing.
- Leads to similar behavior as CountMin sketches, on overloading.

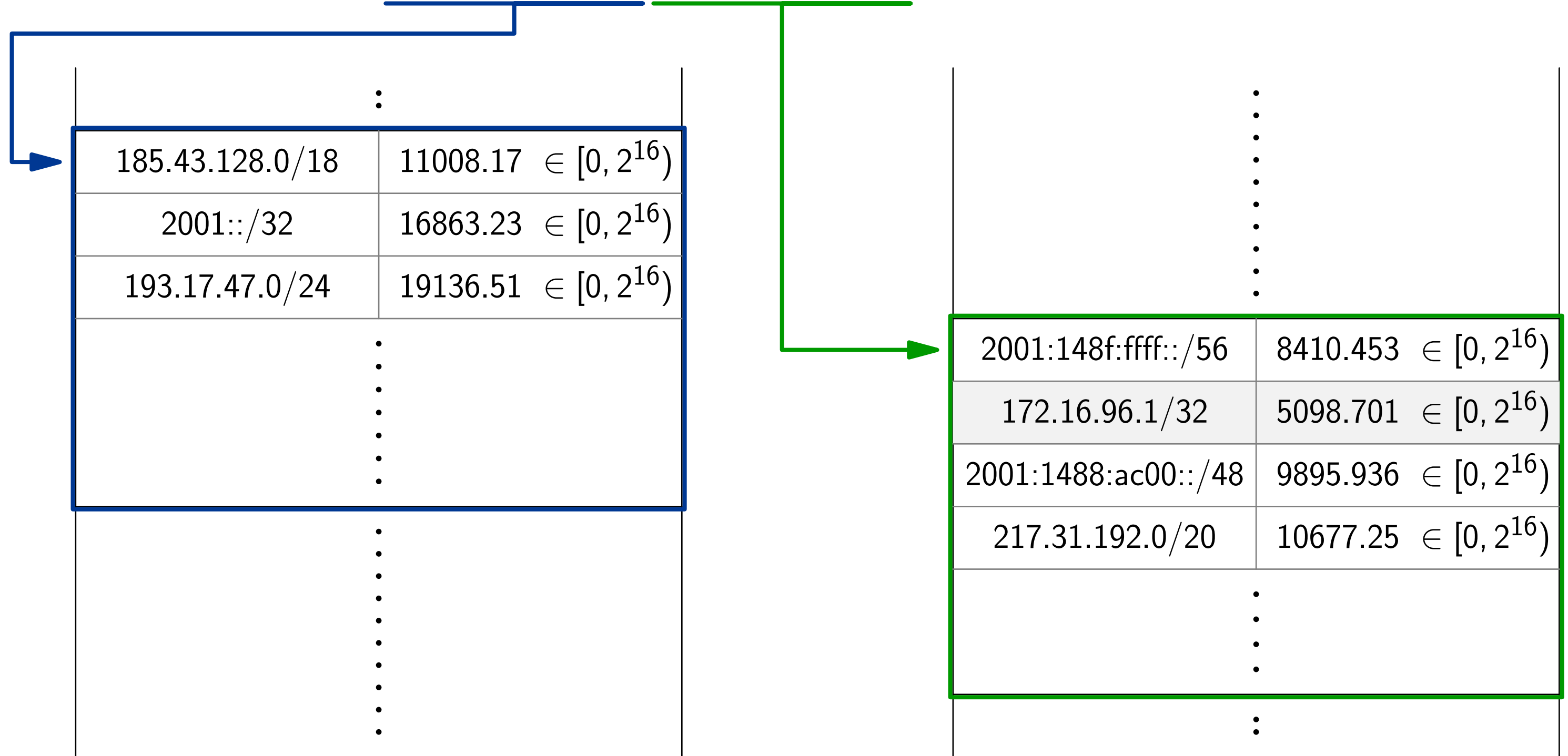
- hashing
 - buckets
 - two tables
- evicting
 - normalized limits
 - choosing minimal
 - **keeping value**

hash(172.16.96.1/32) = 10111101100001101010110000110101001110100111010011101101001000001011111



Implementation

hash(172.16.96.1/32) = 1011110110000110101011000011010100111010011101101001000001011111



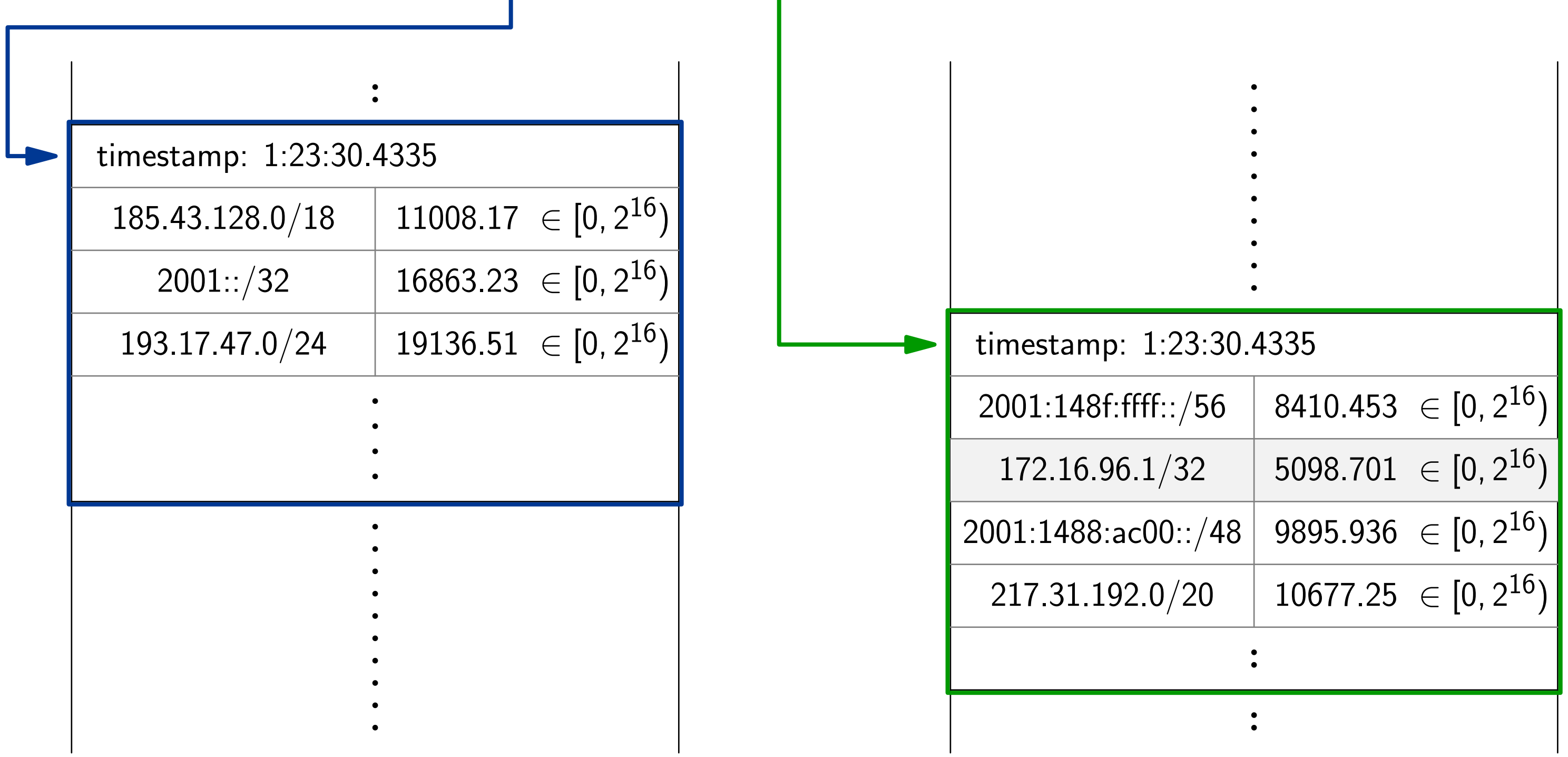
- hashing
 - buckets
 - two tables
- evicting
 - normalized limits
 - choosing minimal
 - keeping value
- lazy decay

- Store timestamp of last decay in each bucket.
- Perform decay on all bucket items.

Implementation

- hashing
 - buckets
 - two tables
- evicting
 - normalized limits
 - choosing minimal
 - keeping value
- lazy decay

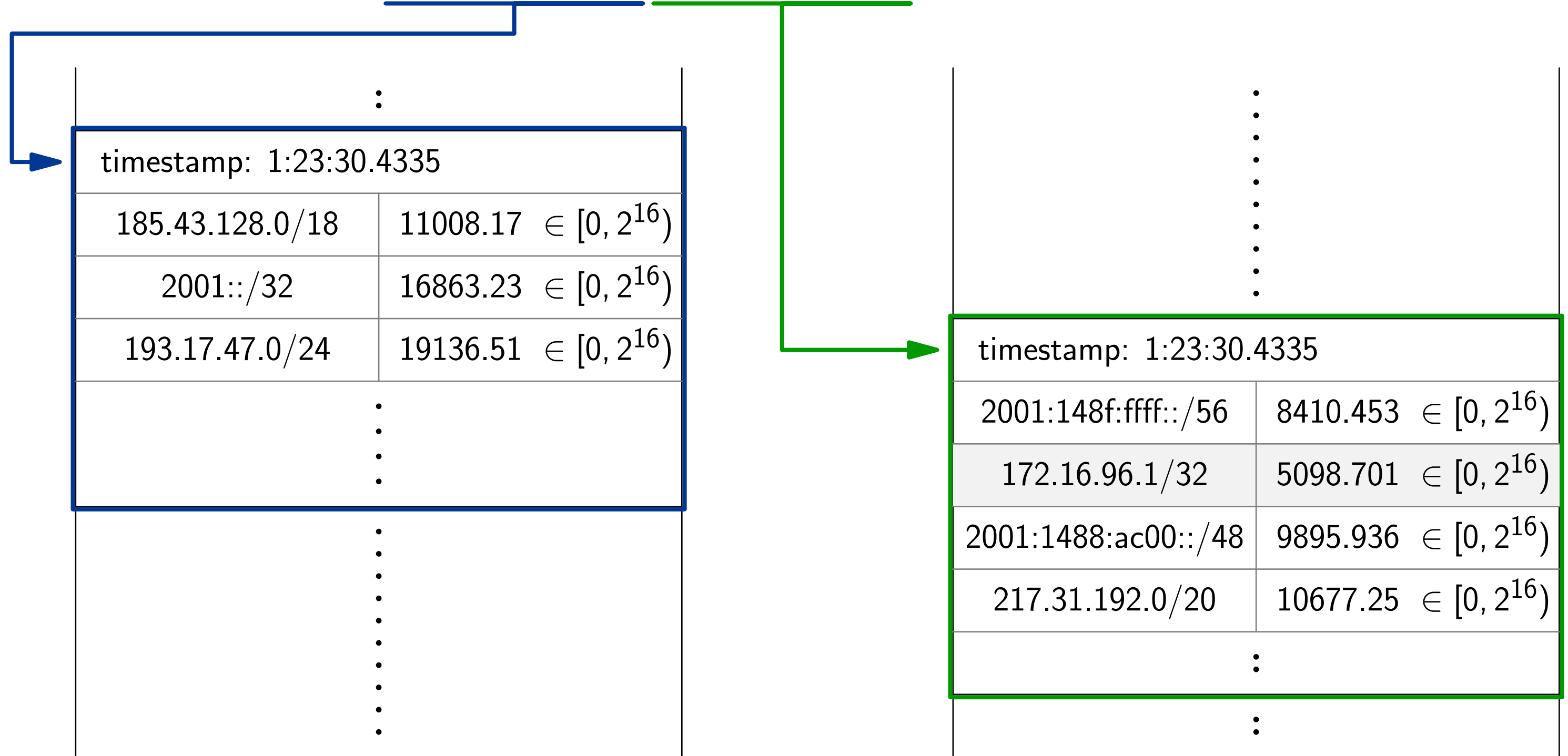
hash(172.16.96.1/32) = 1011110110000110101011000011010100111010011101101001000001011111



Implementation

- hashing
 - buckets
 - two tables
- evicting
 - normalized limits
 - choosing minimal
 - keeping value
- lazy decay
- memory layout

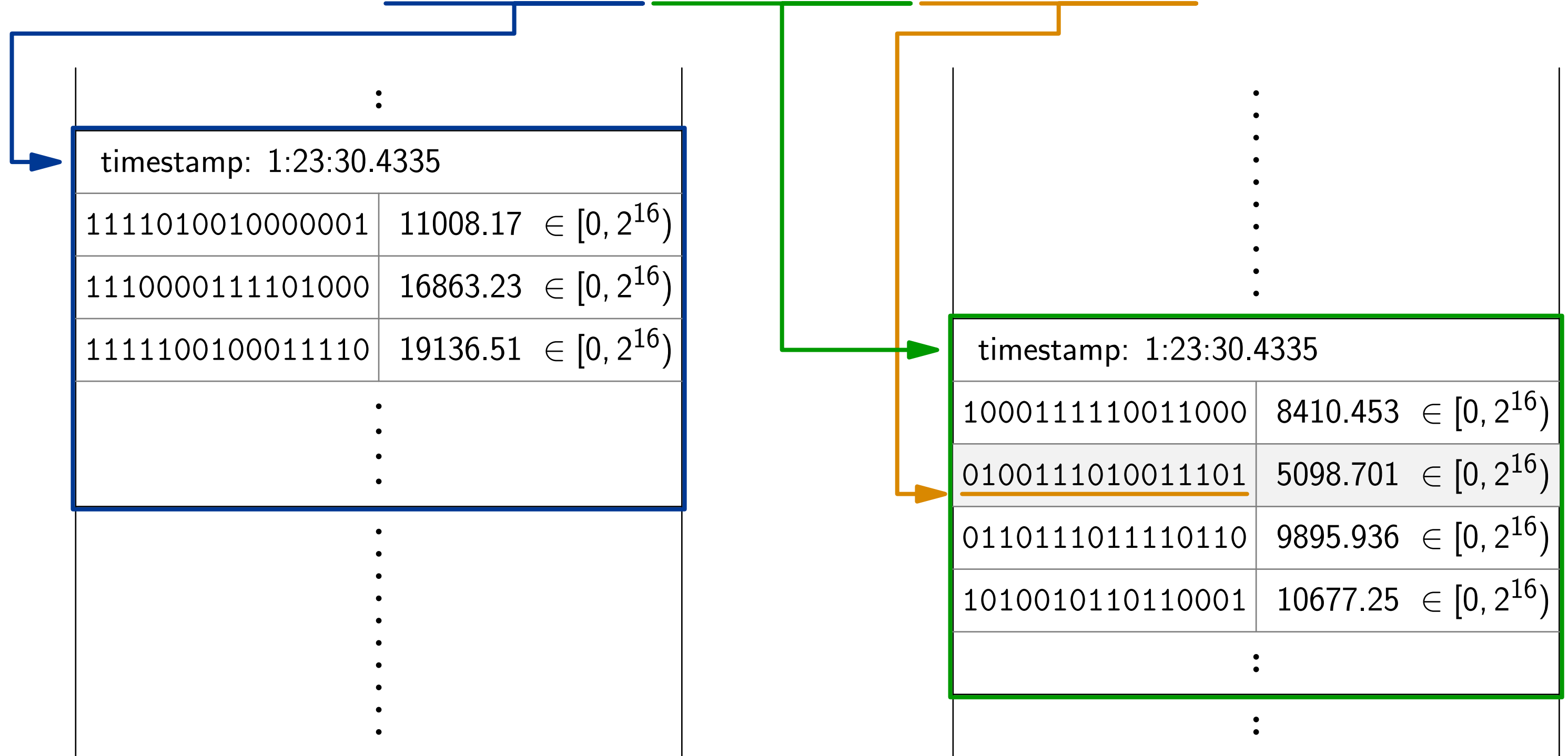
hash(172.16.96.1/32) = 1011110110000110101011000011010100111010011101101001000001011111



Implementation

- Addresses too long, store just another part of their hash (16 bits).
- Collisions may cause sharing counters, but they are very infrequent.

hash(172.16.96.1/32) = 1011110110000110101011000011010100111010011101101001000001011111



- hashing
 - buckets
 - two tables
- evicting
 - normalized limits
 - choosing minimal
 - keeping value
- lazy decay
- memory layout
 - **hashed labels**

Implementation

- Store just the whole part of values, > but rounded probabilistically by the fractional part.
- All calculations 32-bit – 16 in integer part, 16 in fractional.
- Before storing the integer part, use the fractional part as probability of rounding up.
 - Compare the 16 low significant bits to rnd.

- hashing
 - buckets
 - two tables
- evicting
 - normalized limits
 - choosing minimal
 - keeping value
- lazy decay
- memory layout
 - hashed labels
 - **prob. rounding**

hash(172.16.96.1/32) = 101111011000011010101100001101010011

:	
timestamp: 1:23:30.4335	
1111010010000001	11008.17 ∈ [0, 2 ¹⁶)
1110000111101000	16863.23 ∈ [0, 2 ¹⁶)
1111100100011110	19136.51 ∈ [0, 2 ¹⁶)
:	:
:	:
:	:
:	:
:	:
:	:
:	:
:	:

rounding: 5098.7008

16-bit base value 16-bit probability of rounding up (70.08 %)

timestamp: 1:23:30.4335	
1000111110011000	8410.453 ∈ [0, 2 ¹⁶)
<u>0100111010011101</u>	5098.701 ∈ [0, 2 ¹⁶)
0110111011110110	9895.936 ∈ [0, 2 ¹⁶)
1010010110110001	10677.25 ∈ [0, 2 ¹⁶)
:	:
:	:

Implementation

- So we have 16-bit values, but can increment even by much smaller fractions.
- Still very precise – 2^{16} ones required to perform limiting.

- hashing
 - buckets
 - two tables
- evicting
 - normalized limits
 - choosing minimal
 - keeping value
- lazy decay
- memory layout
 - hashed labels
 - **prob. rounding**

hash(172.16.96.1/32) = 101111011000011010101100001101010011

:		
timestamp: 1:23:30.4335		
1111010010000001	[11008.17]	< 2^{16}
1110000111101000	[16863.23]	< 2^{16}
1111100100011110	[19136.51]	< 2^{16}
:	:	:
:	:	:
:	:	:
:	:	:
:	:	:
:	:	:
:	:	:

rounding: 5098.7008

16-bit base value 16-bit probability of rounding up (70.08 %)

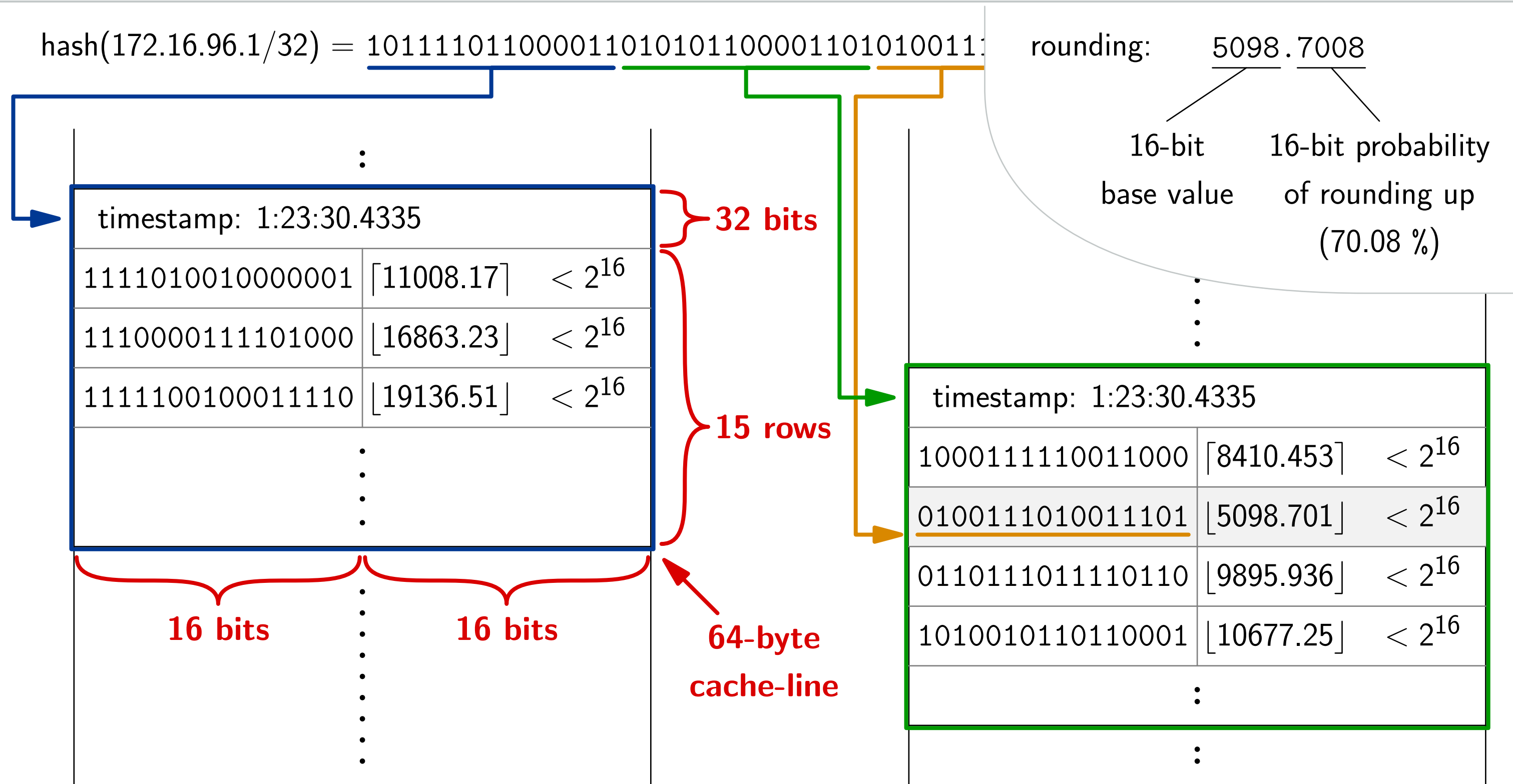
⋮

timestamp: 1:23:30.4335		
1000111110011000	[8410.453]	< 2^{16}
<u>0100111010011101</u>	[5098.701]	< 2^{16}
0110111011110110	[9895.936]	< 2^{16}
1010010110110001	[10677.25]	< 2^{16}
:	:	:
:	:	:

Implementation

- 15x 16-bit label and 16-bit value + 32-bit timestamp is just cache-line.
- Just 2 cache-lines per prefix needed for request, at most 10 in total.

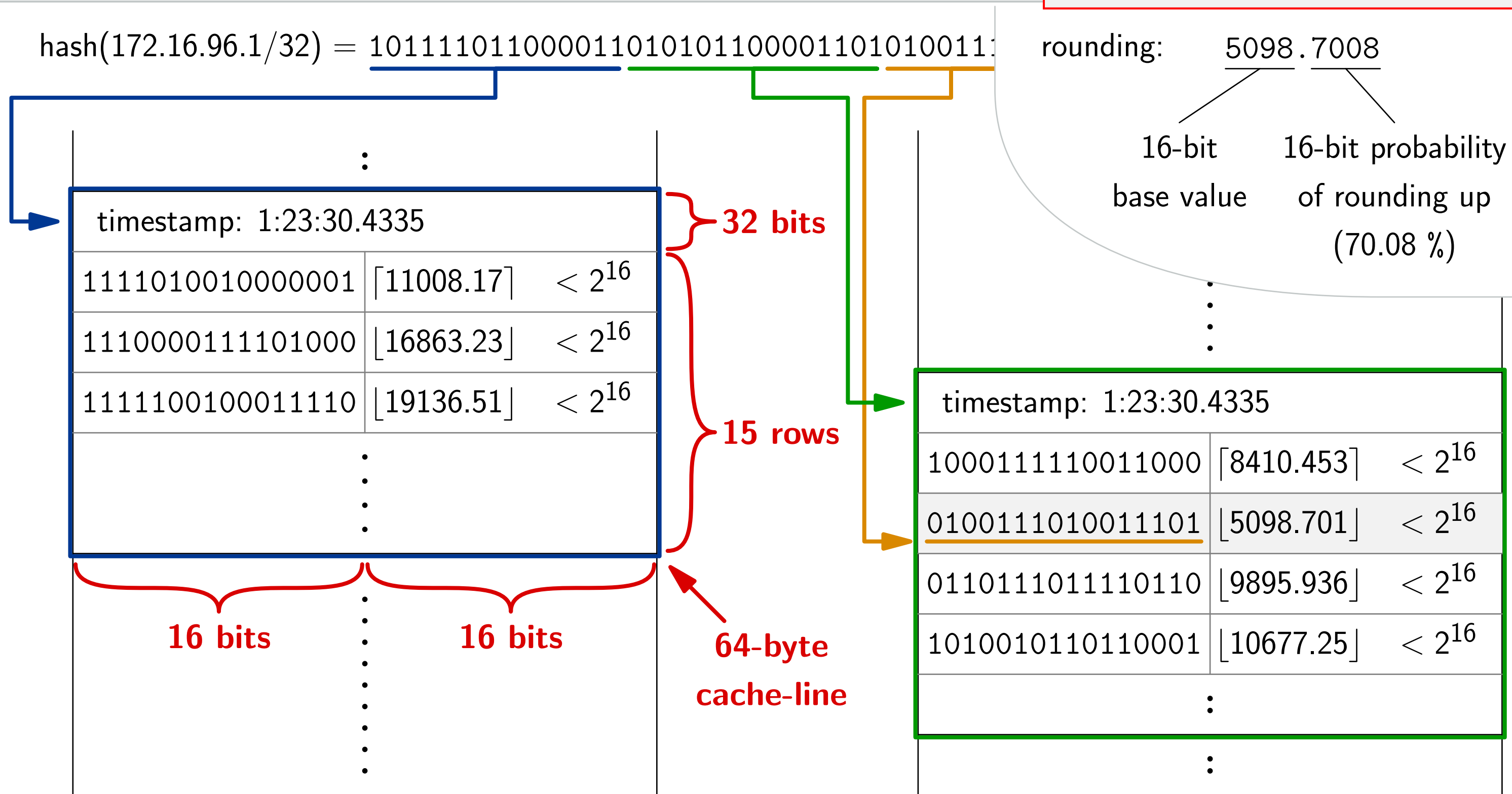
- hashing
 - buckets
 - two tables
- evicting
 - normalized limits
 - choosing minimal
 - keeping value
- lazy decay
- memory layout
 - hashed labels
 - prob. rounding
 - **fit in cache-line**



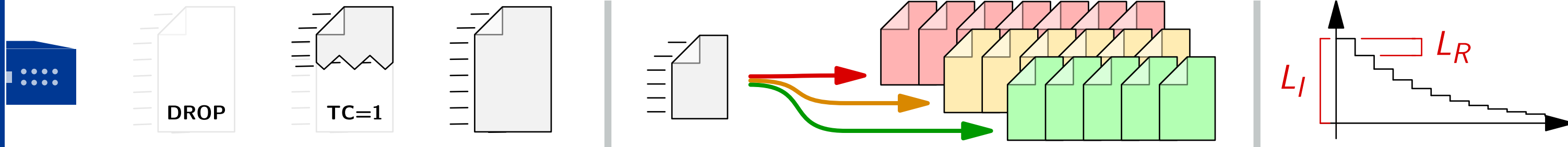
Implementation

- Further optimizations include:
 - Prefetching the 10 cache-lines before accessing.
 - Using atomic instructions to make the data structure lock-free.
 - Using vector instructions to speedup searching in buckets.

- hashing
 - buckets
 - two tables
- evicting
 - normalized limits
 - choosing minimal
 - keeping value
- lazy decay
- memory layout
 - hashed labels
 - prob. rounding
 - fit in cache-line
- optimizations
 - prefetching
 - lock-free
 - vectorization



Summary



- rate-limiting
 - counting UDP queries
 - truncating or dropping
- prioritization
 - measuring time
 - reordering
- counters
 - instant/rate limit
 - exponential decay
 - higher limits for shorter prefixes
- implementation ⇒

hash(172.16.96.1/32) = 101111011000011010101100001101010011

rounding: $\frac{5098.7008}{16\text{-bit base value}}$ $\frac{16\text{-bit probability of rounding up (70.08\%)}$

timestamp: 1:23:30.4335		
1111010010000001	[11008.17]	$< 2^{16}$
1110000111101000	[16863.23]	$< 2^{16}$
1111100100011110	[19136.51]	$< 2^{16}$
⋮	⋮	⋮

32 bits

15 rows

16 bits

16 bits

64-byte cache-line

timestamp: 1:23:30.4335		
1000111110011000	[8410.453]	$< 2^{16}$
<u>0100111010011101</u>	[5098.701]	$< 2^{16}$
0110111011110110	[9895.936]	$< 2^{16}$
1010010110110001	[10677.25]	$< 2^{16}$
⋮	⋮	⋮