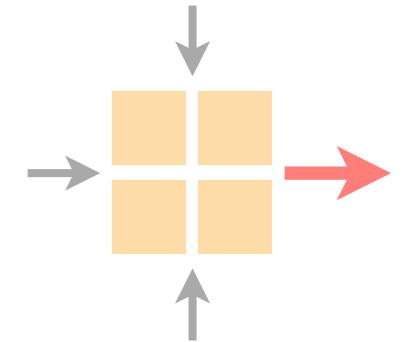


Advanced Topics in Communication Networks



Laurent Vanbever
nsg.ee.ethz.ch

ETH Zürich
Tue 11 Oct 2022

Last week on
Advanced Topics in Communication Networks

We finished our exploration of the P4 language
then looked at our first technique: load balancing

stateful
objects

load balancing

How do we maintain
state in P4?

How do we balance
traffic in a network?

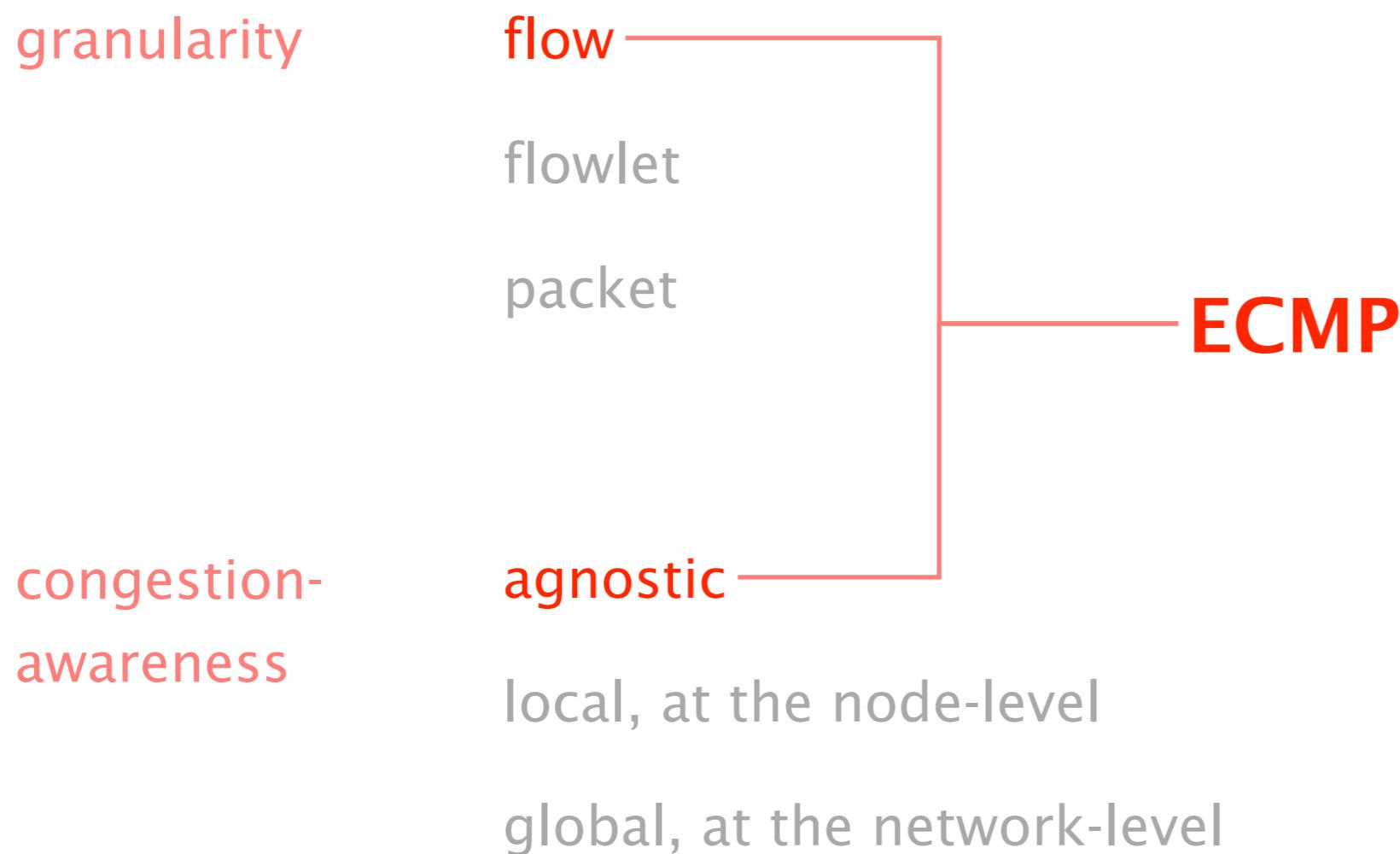
stateful
objects

How do we maintain
state in P4?

load balancing

How do we balance
traffic in a network?

ECMP load-balancing decisions are suboptimal because they are too coarse-grained and congestion-agnostic



load-balancing decision...

granularity	flow
	flowlet
	packet
congestion-awareness	agnostic
	local , at the node-level
	global , at the network-level

load-balancing decision...

granularity

~~flow~~

flowlet

packet ————— finest granularity... but
bad for TCP because of reordering

congestion-
awareness

~~agnostic~~

local, at the node-level

global, at the network-level

load-balancing decision...

granularity	flow
	flowlet
	packet
congestion-awareness	agnostic
	local , at the node-level
	global , at the network-level

load-balancing decision...

granularity

~~flow~~

flowlet

~~packet~~

congestion-awareness

~~agnostic~~

~~local, at the node-level~~

global, at the network-level

handling asymmetry mandates
non-local knowledge

load-balancing decision...

granularity

~~flow~~

flowlet

best pick!

~~packet~~

congestion-awareness

~~agnostic~~

~~local, at the node-level~~

global

best pick!

We started to look at how to implement flowlet-based, congestion-aware load balancing

CONGA: Distributed Congestion-Aware Load Balancing for Datacenters

Mohammad Alizadeh, Tom Edsall, Sarang Dharmapurikar, Ramanan Vaideyanathan, Kevin Chu, Andy Fingerhut, Vinh The Lam (Google), Francis Matus, Rong Pan, Navindra Yadav, George Varghese (Microsoft)

Cisco Systems

ABSTRACT

We present the design, implementation, and evaluation of CONGA, a network-based distributed congestion-aware load balancing mechanism for datacenters. CONGA exploits recent trends including the use of regular Clos topologies and overlays for network virtualization. It splits TCP flows into flowlets, estimates real-time congestion on fabric paths, and allocates flowlets to paths based on feedback from remote switches. This enables CONGA to efficiently balance load and seamlessly handle asymmetry, without requiring any TCP modifications. CONGA has been implemented in custom ASICs and tested in a new datacenter fabric described in experiments. CONGA has 8x better flow completion times than ECMP even with a single link failure and achieves 2-8x better throughput than MPTCP in bursty scenarios. Further, the Price of Anarchy for CONGA is provably small in Leaf-Spine topologies; hence CONGA is nearly as effective as a centralized scheduler while being able to react to congestion in microseconds. Our main thesis is that datacenter fabric load balancing is best done in the network, and requires global schemes such as CONGA to handle asymmetry.

Categories and Subject Descriptors: C.2.1 [Computer-Communication Networks]: Network Architecture and Design
Keywords: Datacenter fabric; Load balancing; Distributed

1. INTRODUCTION

Datacenter networks being deployed by cloud providers as well as enterprises must provide large bisection bandwidth to support an ever-increasing number of applications from financial services to big-data analytics. They also must provide agility, enabling any application to be deployed at any server, in order to realize operational efficiency and reduce costs. Seminal papers such as VL2 [18] and Portland [1] showed how to achieve this with Clos topologies. Equal Cost MultiPath (ECMP) load balancing, and the decoupling of endpoint addresses from their location. These design principles are followed by next generation overlay technologies that accomplish the same goals using standard encapsulations such as VXLAN [35] and NVGRE [45].

However, it is well known [2, 41, 9, 27, 44, 10] that ECMP can balance load poorly. First, because ECMP randomly hashes flows to paths, hash collisions can cause significant imbalance if there are a few large flows. More importantly, ECMP uses a purely *local* decision to split traffic among local network paths without knowledge of potential downstream congestion on each path. Thus ECMP fares poorly with *asymmetry* caused by link failures that occur frequently and are disruptive in datacenters [17, 34]. For instance, the recent study by Gill *et al.* [17] shows that failures can reduce delivered traffic by up to 40% despite built-in redundancy.

Broadly speaking, the prior work on addressing ECMP's shortcomings can be classified as either centralized scheduling (e.g., Hedera [2]), local switch mechanisms (e.g., Flare [27]), or host-based transport protocols (e.g., MPTCP [41]). These approaches all have important drawbacks. Centralized schemes are too slow for the traffic volatility in datacenters [28, 8] and local congestion-aware mechanisms are suboptimal and can perform even worse than ECMP with asymmetry (§2.4). Host-based methods such as MPTCP are challenging to deploy because network operators often do not control the end-host stack (e.g., in a public cloud) and even when they do, some high performance applications (such as low latency storage systems [39, 7]) bypass the kernel and implement their own transport. Further, host-based load balancing adds more complexity to an already complex transport layer burdened by new requirements such as low latency and burst tolerance [4] in datacenters. As our experiments with MPTCP show, this can make for brittle and unreliable schemes [5].

Thus from a philosophical standpoint it is worth asking: Can load balancing be done in the network without adding to the complexity of the transport layer? Can such a network-based approach compute globally optimal allocations, and yet be implementable in a realizable and distributed fashion to allow rapid reaction in microseconds? Can such a mechanism be deployed today using standard encapsulation formats? We seek to answer these questions in this paper with a new scheme called CONGA (for Congestion Aware Balancing). CONGA has been implemented in custom ASICs for a major new datacenter fabric product line. While we report on lab experiments using working hardware together with simulations and mathematical analysis, customer trials are scheduled in a few months as of the time of this writing.

Figure 1 surveys the design space for load balancing and places CONGA in context by following the thick red lines through the design tree. At the highest level, CONGA is a distributed scheme to allow rapid round-trip timescale reaction to congestion to cope with bursty datacenter traffic [28, 8]. CONGA is implemented within the network to avoid the deployment issues of host-based methods and additional complexity in the transport layer. To deal with asymmetry, unlike earlier proposals such as Flare [27] and LocalFlow [44] that only use local information, CONGA uses global congestion information, a design choice justified in detail in §2.4.

Permission to make digital or hard copies of all or part of this work for personal classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyright for works published in this work owned by others than ACM may be reserved or transferred only with explicit permission or a fee. To copy otherwise, or to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from permissions@acm.org.
SIGCOMM'14, August 17–22, 2014, Chicago, IL, USA
Copyright 2014 ACM 978-1-4503-2836-4/14/08 .\\$15.00.
http://dx.doi.org/10.1145/2619259.2626116.

Let it Flow: Resilient Asymmetric Load Balancing with Flowlet Switching

Erico Vanini*, Rong Pan*, Mohammad Alizadeh†, Parvin Taheri*, Tom Edsall*

*Cisco Systems †Massachusetts Institute of Technology

Abstract

Datacenter networks require efficient multi-path load balancing to achieve high bisection bandwidth. Despite much progress in recent years towards addressing this challenge, a load balancing design that is both simple to implement and resilient to network asymmetry has remained elusive. In this paper, we show that *flowlet switching*, an idea first proposed more than a decade ago, is a powerful technique for resilient load balancing with asymmetry. Flowlets have a remarkable *elasticity* property: their size changes automatically based on traffic conditions on their path. We use this insight to develop LetFlow, a very simple load balancing scheme that is resilient to asymmetry. LetFlow simply picks paths at random for flowlets and lets their elasticity naturally balance the traffic on different paths. Our extensive evaluation with real hardware and packet-level simulations shows that LetFlow is very effective. Despite being much simpler, it performs significantly better than other traffic oblivious schemes like WCMP and Presto in asymmetric scenarios, while achieving average flow completions time within 10-20% of CONGA in testbed experiments and 2x of CONGA in simulated topologies with large asymmetry and heavy traffic load.

1 Introduction

Datacenter networks must provide large bisection bandwidth to support the increasing traffic demands of applications such as big-data analytics, web services, and cloud storage. They achieve this by load balancing traffic over many paths in multi-rooted tree topologies such as Clos [13] and Fat-tree [1]. These designs are widely deployed; for instance, Google has reported on using Clos fabrics with more than 1 Pbps of bisection bandwidth in its datacenters [25].

The standard load balancing scheme in today's datacenters, Equal Cost MultiPath (ECMP) [16], randomly assigns flows to different paths using a hash taken over packet headers. ECMP is widely deployed due to its simplicity but suffers from well-known performance problems such as hash collisions and the inability to adapt to asymmetry in the network topology. A rich body of work [10, 2, 22, 23, 18, 3, 15, 21] has thus emerged on

better load balancing designs for datacenter networks.

A defining feature of these designs is the information that they use to make decisions. At one end of the spectrum are designs that are oblivious to traffic conditions [16, 10, 9, 15] or rely only on local measurements [24, 20] at the switches. ECMP and Presto [15], which picks paths in round-robin fashion for fixed-sized chunks of data (called "flowcells"), fall in this category. At the other extreme are designs [2, 22, 23, 18, 3, 21, 29] that use knowledge of traffic conditions and congestion on different paths to make decisions. Two recent examples are CONGA [3] and HULA [21], which use feedback between the switches to gather path-wise congestion information and shift traffic to less-congested paths.

Load balancing schemes that require path congestion information, naturally, are much more complex. Current designs either use a centralized fabric controller [2, 8, 22] to optimize path choices frequently or require non-trivial mechanisms, at the end-hosts [23, 18] or switches [3, 21, 30], to implement end-to-end or hop-by-hop feedback. On the other hand, schemes that lack visibility into path congestion have a key drawback: they perform poorly in *asymmetric topologies* [3]. As we discuss in §2.1, the reason is that the optimal traffic split across asymmetric paths depends on (dynamically varying) traffic conditions; hence, traffic-oblivious schemes are fundamentally unable to make optimal decisions and can perform poorly in asymmetric topologies.

Asymmetry is common in practice for a variety of reasons, such as link failures and heterogeneity in network equipment [31, 12, 3]. Handling asymmetry gracefully, therefore, is important. This raises the question: *are there simple load balancing schemes that are resilient to asymmetry?* In this paper, we answer this question in the affirmative by developing LetFlow, a simple scheme that requires no state to make load balancing decisions and yet it is very resilient to network asymmetry.

LetFlow is extremely simple: switches pick a path at random for each *flowlet*. That's it! A flowlet is a burst of packets that is separated in time from other bursts by a sufficient gap — called the “flowlet timeout”. Flowlet switching [27, 20] was proposed over a decade ago as a way to split TCP flows across multiple paths without causing packet reordering. Remarkably, as we uncover in this paper, flowlet switching is also a powerful technique

CONGA [SIGCOMM'14]

LetFlow [NSDI'17]

CONGA: Distributed Congestion-Aware Load Balancing for Datacenters

Mohammad Alizadeh, Tom Edsall, Sarang Dharmapurikar, Ramanan Vaideyanathan, Kevin Chu, Andy Fingerhut, Vinh The Lam (Google), Francis Matus, Rong Pan, Navindra Yadav, George Varghese (Microsoft)

Cisco Systems

ABSTRACT

We present the design, implementation, and evaluation of CONGA, a network-based distributed congestion-aware load balancing mechanism for datacenters. CONGA exploits recent trends including the use of regular Clos topologies and overlays for network virtualization. It splits TCP flows into flowlets, estimates real-time congestion on fabric paths, and allocates flowlets to paths based on feedback from remote switches. This enables CONGA to efficiently balance load and seamlessly handle asymmetry, without requiring any TCP modifications. CONGA has been implemented in custom ASICs and is being deployed in several described experiments. CONGA has 8x better flow completion times than ECMP even with a single link failure and achieves 2-8x better throughput than MPTCP in bursty scenarios. Further, the Price of Anarchy for CONGA is provably small in Leaf-Spine topologies; hence CONGA is nearly as effective as a centralized scheduler while being able to react to congestion in microseconds. Our main thesis is that datacenter fabric load balancing is best done in the network, and requires global schemes such as CONGA to handle asymmetry.

Categories and Subject Descriptors: C.2.1 [Computer-Communication Networks]: Network Architecture and Design
Keywords: Datacenter fabric; Load balancing; Distributed

1. INTRODUCTION

Datacenter networks being deployed by cloud providers as well as enterprises must provide large bisection bandwidth to support an ever-increasing number of applications from financial trading to big-data analytics. They also must provide agility, enabling any application to be deployed at any server, in order to realize operational efficiency and reduce costs. Seminal papers such as VL2 [18] and Portland [1] showed how to achieve this with Clos topologies. Equal Cost MultiPath (ECMP) load balancing, and the decoupling of endpoint addresses from their location. These design principles are followed by next generation overlay technologies that accomplish the same goals using standard encapsulations such as VXLAN [35] and NVGRE [45].

However, it is well known [2, 41, 9, 27, 44, 10] that ECMP can balance load poorly. First, because ECMP randomly hashes flows to paths, hash collisions can cause significant imbalance if there are a few large flows. More importantly, ECMP uses a purely *local* decision to split traffic among local paths without knowledge of potential downstream congestion on each path. Thus ECMP fares poorly with *asymmetry* caused by link failures that occur frequently and are disruptive in datacenters [17, 34]. For instance, the recent study by Gill *et al.* [17] shows that failures can reduce delivered traffic by up to 40% despite built-in redundancy.

Broadly speaking, the prior work on addressing ECMP's shortcomings can be classified as either centralized scheduling (e.g., Hedera [2]), local switch mechanisms (e.g., Flare [27]), or host-based transport protocols (e.g., MPTCP [41]). These approaches all have important drawbacks. Centralized schemes are too slow for the traffic volatility in datacenters [28, 8] and local congestion-aware mechanisms are suboptimal and can perform even worse than ECMP with asymmetry (§2.4). Host-based methods such as MPTCP are challenging to deploy because network operators often do not control the end-host stack (e.g., in a public cloud) and even when they do, some high performance applications (such as low latency storage systems [39, 7]) bypass the kernel and implement their own transport. Further, host-based load balancing adds more complexity to an already complex transport layer burdened by new requirements such as low latency and burst tolerance [4] in datacenters. As our experiments with MPTCP show, this can make for brittle and unreliable designs (§2.5).

Thus from a philosophical standpoint it is worth asking: Can load balancing be done in the network without adding to the complexity of the transport layer? Can such a network-based approach compute globally optimal allocations, and yet be implementable in a realizable and distributed fashion to allow rapid reaction in microseconds? Can such a mechanism be deployed today using standard encapsulation formats? We seek to answer these questions in this paper with a new scheme called CONGA (for Congestion Aware Balancing). CONGA has been implemented in custom ASICs for a major new datacenter fabric product line. While we report on lab experiments using working hardware together with simulations and mathematical analysis, customer trials are scheduled in a few months as of the time of this writing.

Figure 1 surveys the design space for load balancing and places CONGA in context by following the thick red lines through the design tree. At the highest level, CONGA is a distributed scheme to allow rapid round-trip timescale reaction to congestion to cope with bursty datacenter traffic [28, 8]. CONGA is implemented within the network to avoid the deployment issues of host-based methods and additional complexity in the transport layer. To deal with asymmetry, unlike earlier proposals such as Flare [27] and LocalFlow [44] that only use local information, CONGA uses global congestion information, a design choice justified in detail in §2.4.

Permission to make digital or hard copies of all or part of this work for personal classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyright for works published in this work owned by others than ACM may be reserved or transferred only with explicit permission or written re-publication, to post on servers or redistribute to lists, requires prior specific permission and/or a fee. Request permissions from permissions@acm.org.
SIGCOMM'14, August 17–22, 2014, Chicago, IL, USA
Copyright 2014 ACM 978-1-4503-2836-4/14/08 . \$15.00.
http://dx.doi.org/10.1145/2619259.2626316.

Let it Flow: Resilient Asymmetric Load Balancing with Flowlet Switching

Erico Vanini*, Rong Pan*, Mohammad Alizadeh†, Parvin Taheri*, Tom Edsall*

*Cisco Systems †Massachusetts Institute of Technology

Abstract

Datacenter networks require efficient multi-path load balancing to achieve high bisection bandwidth. Despite much progress in recent years towards addressing this challenge, a load balancing design that is both simple to implement and resilient to network asymmetry has remained elusive. In this paper, we show that *flowlet switching*, an idea first proposed more than a decade ago, is a powerful technique for resilient load balancing with asymmetry. Flowlets have a remarkable *elasticity* property: their size changes automatically based on traffic conditions on their path. We use this insight to develop LetFlow, a very simple load balancing scheme that is resilient to asymmetry. LetFlow simply picks paths at random for flowlets and lets their elasticity naturally balance the traffic on different paths. Our extensive evaluation with real hardware and packet-level simulations shows that LetFlow is very effective. Despite being much simpler, it performs significantly better than other traffic oblivious schemes like WCMP and Presto in asymmetric scenarios, while achieving average flow completions time within 10-20% of CONGA in tested experiments and 2x of CONGA in simulated topologies with large asymmetry and heavy traffic load.

1 Introduction

Datacenter networks must provide large bisection bandwidth to support the increasing traffic demands of applications such as big-data analytics, web services, and cloud storage. They achieve this by load balancing traffic over many paths in multi-rooted tree topologies such as Clos [13] and Fat-tree [1]. These designs are widely deployed; for instance, Google has reported on using Clos fabrics with more than 1 Pbps of bisection bandwidth in its datacenters [25].

The standard load balancing scheme in today's datacenters, Equal Cost MultiPath (ECMP) [16], randomly assigns flows to different paths using a hash taken over packet headers. ECMP is widely deployed due to its simplicity but suffers from well-known performance problems such as hash collisions and the inability to adapt to asymmetry in the network topology. A rich body of work [10, 2, 22, 23, 18, 3, 15, 21] has thus emerged on

better load balancing designs for datacenter networks.

A defining feature of these designs is the information that they use to make decisions. At one end of the spectrum are designs that are oblivious to traffic conditions [16, 10, 9, 15] or rely only on local measurements [24, 20] at the switches. ECMP and Presto [15], which picks paths in round-robin fashion for fixed-sized chunks of data (called "flowcells"), fall in this category. At the other extreme are designs [2, 22, 23, 18, 3, 21, 29] that use knowledge of traffic conditions and congestion on different paths to make decisions. Two recent examples are CONGA [3] and HULA [21], which use feedback between the switches to gather path-wise congestion information and shift traffic to less-congested paths.

Load balancing schemes that require path congestion information, naturally, are much more complex. Current designs either use a centralized fabric controller [2, 8, 22] to optimize path choices frequently or require non-trivial mechanisms, at the end-hosts [23, 18] or switches [3, 21, 30], to implement end-to-end or hop-by-hop feedback. On the other hand, schemes that lack visibility into path congestion have a key drawback: they perform poorly in *asymmetric topologies* [3]. As we discuss in §2.1, the reason is that the optimal traffic split across asymmetric paths depends on (dynamically varying) traffic conditions; hence, traffic-oblivious schemes are fundamentally unable to make optimal decisions and can perform poorly in asymmetric topologies.

Asymmetry is common in practice for a variety of reasons, such as link failures and heterogeneity in network equipment [31, 12, 3]. Handling asymmetry gracefully, therefore, is important. This raises the question: *are there simple load balancing schemes that are resilient to asymmetry?* In this paper, we answer this question in the affirmative by developing LetFlow, a simple scheme that requires no state to make load balancing decisions and yet is very resilient to network asymmetry.

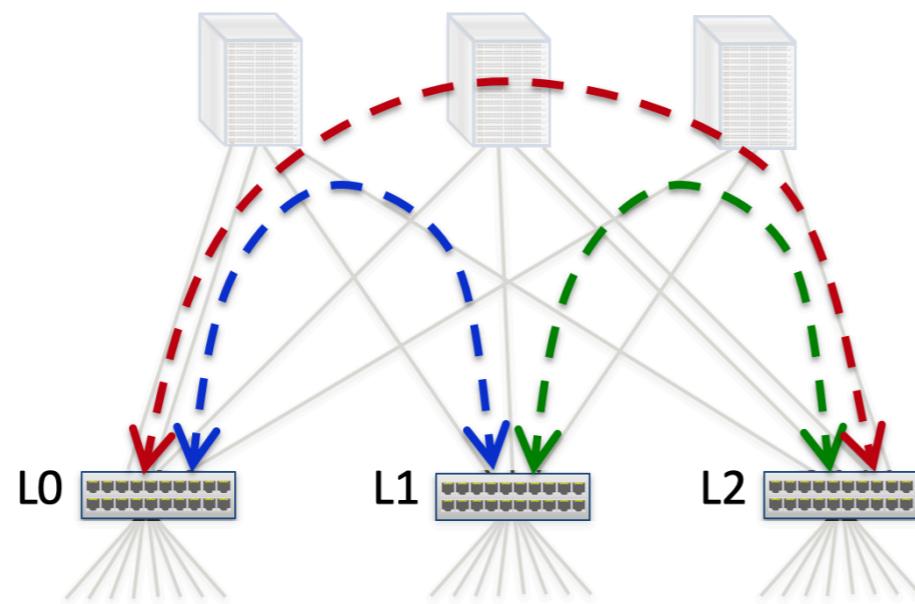
LetFlow is extremely simple: switches pick a path at random for each *flowlet*. That's it! A flowlet is a burst of packets that is separated in time from other bursts by a sufficient gap — called the “flowlet timeout”. Flowlet switching [27, 20] was proposed over a decade ago as a way to split TCP flows across multiple paths without causing packet reordering. Remarkably, as we uncover in this paper, flowlet switching is also a powerful technique

CONGA [SIGCOMM'14]

LetFlow [NSDI'17]

CONGA in 1 Slide

1. Leaf switches (top-of-rack) track congestion to other leaves on different paths **in near real-time**
1. Use greedy decisions to minimize bottleneck util



Fast feedback loops
between leaf switches,
directly in dataplane

12

Source: CONGA: Distributed Congestion-Aware Load Balancing for Datacenters,
Mohammad Alizadeh et al., 2014

This week on
Advanced Topics in Communication Networks

P4 hardware
target

Probabilistic
data structures

P4 hardware
target

Probabilistic
data structures

How do we build a *fast*
reprogrammable switch?

“Programmable switches are 10-100x slower than fixed-function switches. They cost more and consume more power.”

Conventional wisdom in networking

How can we enable network programmability in the field,
at reasonable cost, and without sacrificing speed

supporting Tbps of
backplane throughput

Let's look at a concrete design: Reconfigurable Match Tables (RMT)

The screenshot shows a PDF document titled "Forwarding Metamorphosis: Fast Programmable Match-Action Processing in Hardware for SDN". The document is presented in a window with a toolbar at the top, including icons for zoom, search, and file operations. The title is centered above the author information, which lists Pat Bosshart†, Glen Gibb‡, Hun-Seok Kim†, George Varghese§, Nick McKeown‡, Martin Izzard†, Fernando Mujica†, and Mark Horowitz‡. Below the authors, their institutional affiliations are given: Texas Instruments, Stanford University, and Microsoft Research. The email addresses for the authors are listed as pat.bosshart@gmail.com, {grg, nickm, horowitz}@stanford.edu, varghese@microsoft.com, and {hkim, izzard, fmujica}@ti.com.

ABSTRACT

In Software Defined Networking (SDN) the control plane is physically separate from the forwarding plane. Control software programs the forwarding plane (e.g., switches and routers) using an open interface, such as OpenFlow. This paper aims to overcomes two limitations in current switching chips and the OpenFlow protocol: i) current hardware switches are quite rigid, allowing “Match-Action” processing on only a fixed set of fields, and ii) the OpenFlow specification only defines a limited repertoire of packet processing actions. We propose the RMT (reconfigurable match tables) model, a new RISC-inspired pipelined architecture for switching chips, and we identify the essential minimal set of action primitives to specify how headers are processed in hardware. RMT allows the forwarding plane to be changed in the field without modifying hardware. As in OpenFlow, the programmer can specify multiple match tables of arbitrary width and depth, subject only to an overall resource limit, with each table configurable for matching on arbitrary fields. However, RMT allows the programmer to modify *all* header fields much more comprehensively than in OpenFlow. Our paper describes the design of a 64 port by 10 Gb/s switch chip implementing the RMT model. Our concrete design demonstrates, contrary to concerns within the community, that flexible OpenFlow hardware switch implementations are feasible at almost no additional cost or power.

1. INTRODUCTION

To improve is to change; to be perfect is to change often. — Churchill

Good abstractions—such as virtual memory and time-sharing—are paramount in computer systems because they allow systems to deal with change and allow simplicity of programming at the next higher layer. Networking has progressed because of key abstractions: TCP provides the abstraction of connected queues between endpoints, and IP provides a simple datagram abstraction from an endpoint to the network edge. However, routing and forwarding *within* the network remain a confusing conglomerate of routing protocols (e.g., BGP, ICMP, MPLS) and forwarding behaviors (e.g., routers, bridges, firewalls), and the control and forwarding planes remain intertwined inside closed, vertically integrated boxes.

Software-defined networking (SDN) took a key step in abstracting network functions by separating the roles of the control and forwarding planes via an *open* interface between them (e.g., OpenFlow [27]). The control plane is lifted up and out of the switch, placing it in external software. This programmatic control of the forwarding plane allows network owners to add new functionality to their network, while replicating the behavior of existing protocols. OpenFlow has become quite well-known as an interface between the control plane and the forwarding plane based on the approach known as “Match-Action”. Roughly, a subset of packet bytes

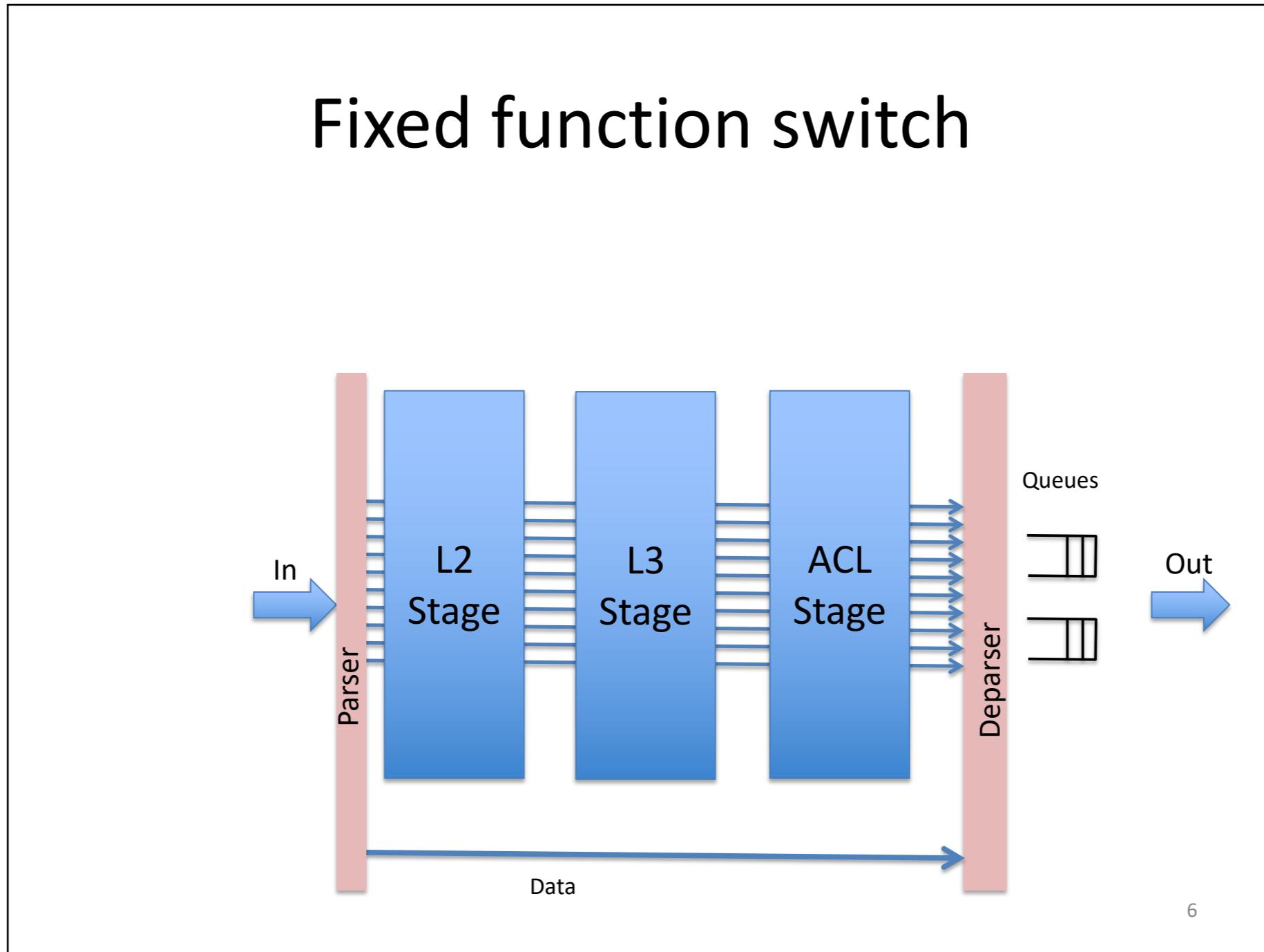
[SIGCOMM'13]

The paper argues that flexibility does **not** come at the price of performance or cost

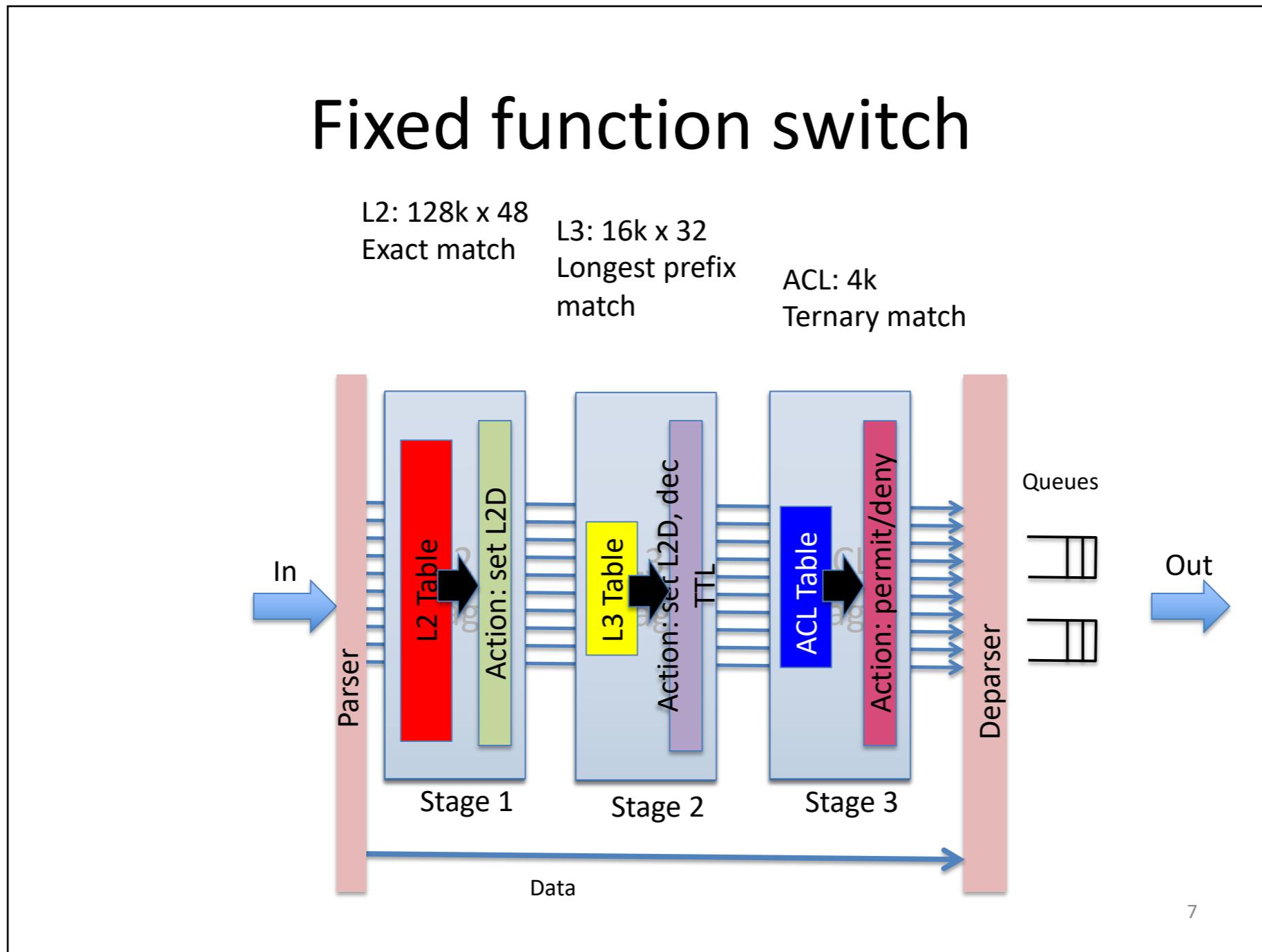
Outline

- Conventional switch chips are inflexible
- SDN demands flexibility...sounds expensive...
- How do we do it: The RMT switch model
- Flexibility costs less than 15%

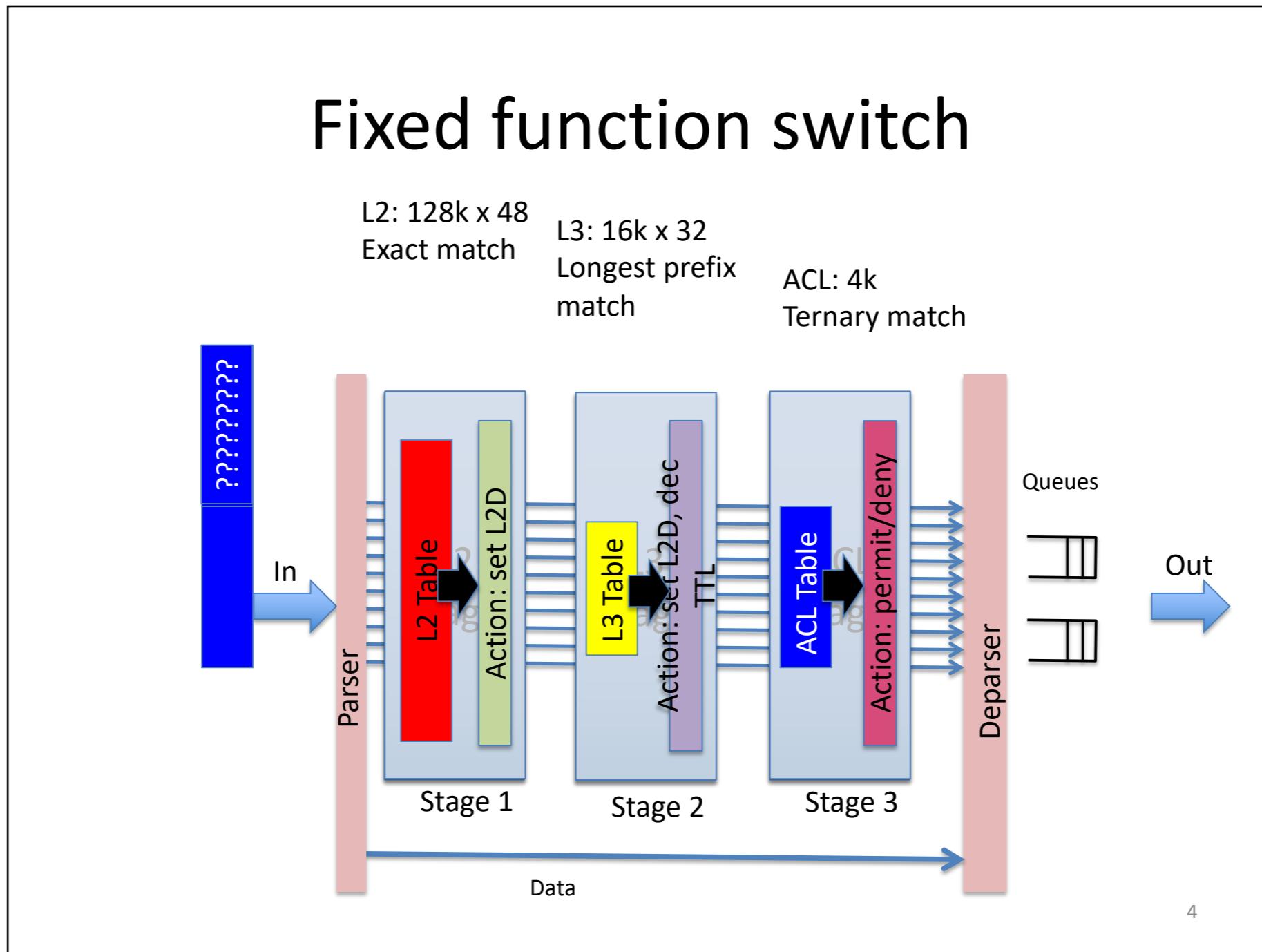
Let's look first at a fixed-function switch composed of a (de-)parser and a sequence of processing stages



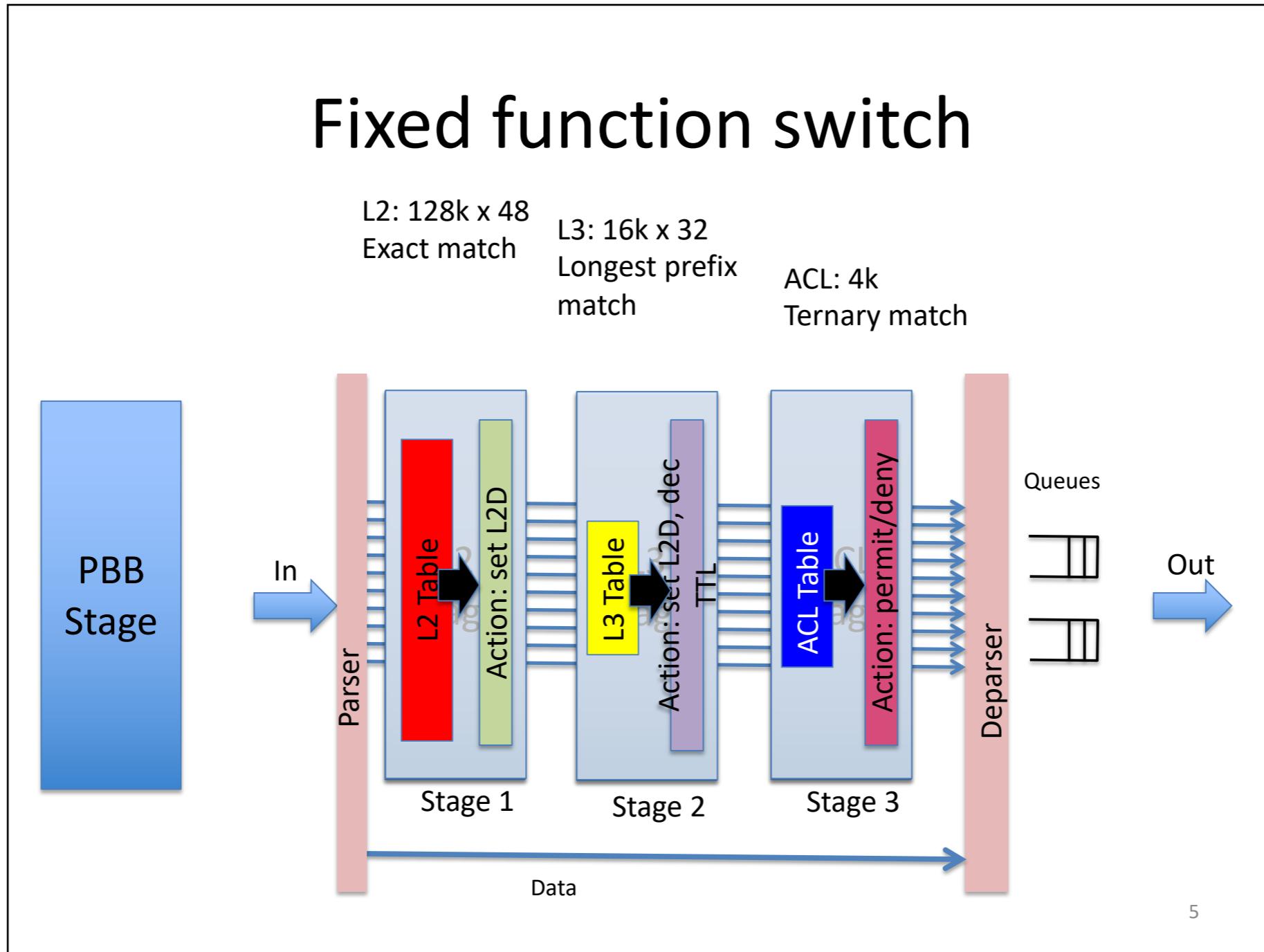
In such a switch,
each stage is particularized to its usage



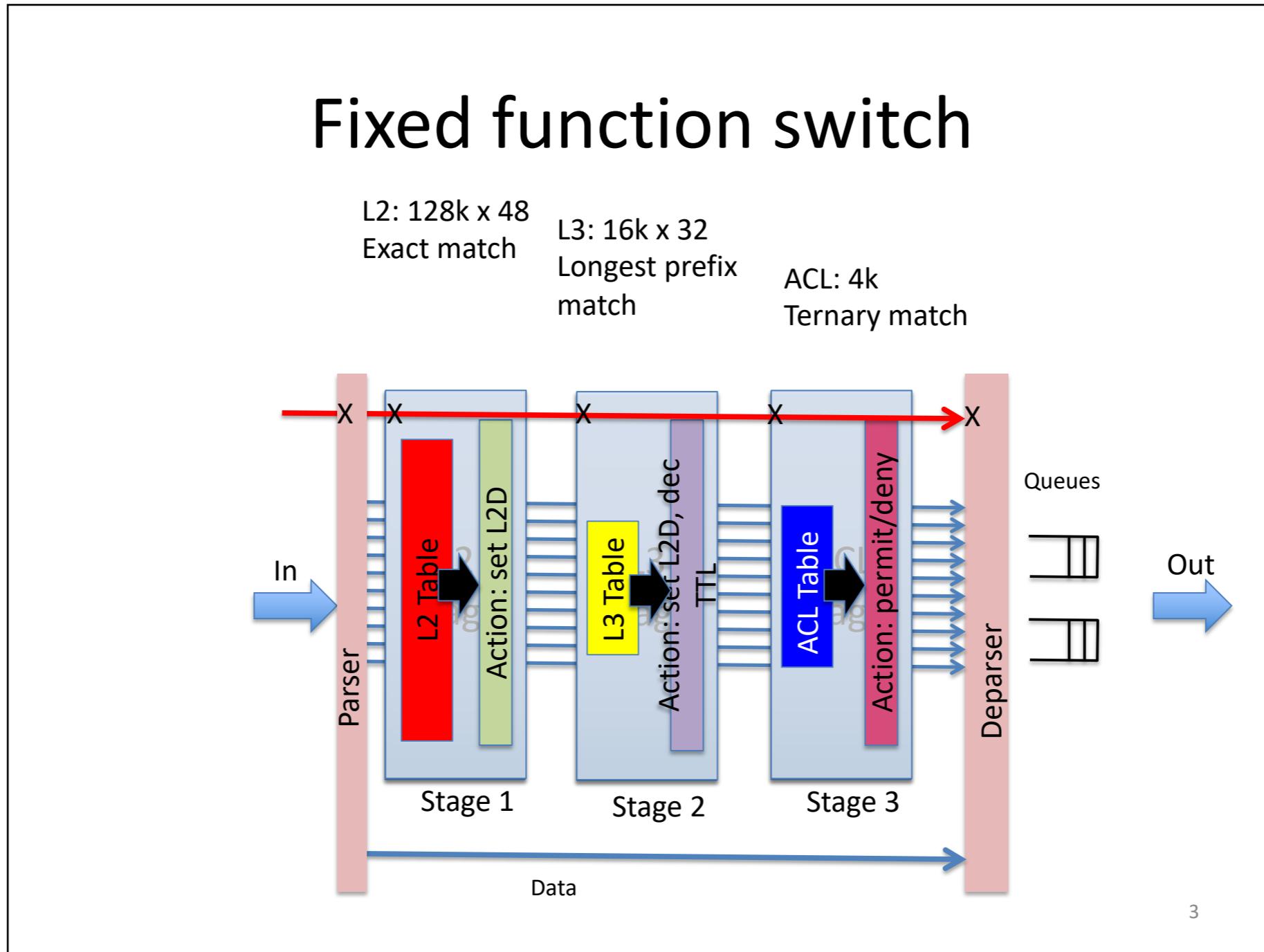
This specificity makes it impossible to...
re-arrange the tables' sizes



This specificity makes it impossible to...
add a new table



This specificity makes it impossible to...
support new headers or new actions



What We Set Out To Learn

- How do I design a flexible switch chip?
- What does the flexibility cost?

Alternative ways to enable flexibility don't compare in terms of cost-performance ratio

What about Alternatives?
Aren't there other ways to get flexibility?

- Software? 100x too slow, expensive
- NPUs? 10x too slow, expensive
- FPGAs? 10x too slow, expensive

Unsurprisingly...
building flexible switching chipset *is* challenging

What's Hard about a Flexible Switch Chip?

- Big chip
- High frequency
- Wiring intensive
- Many crossbars
- Lots of TCAM
- Interaction between physical design and architecture
- Good news? No need to read 7000 IETF RFC's!

Enter...

Reconfigurable Match Tables (RMT)

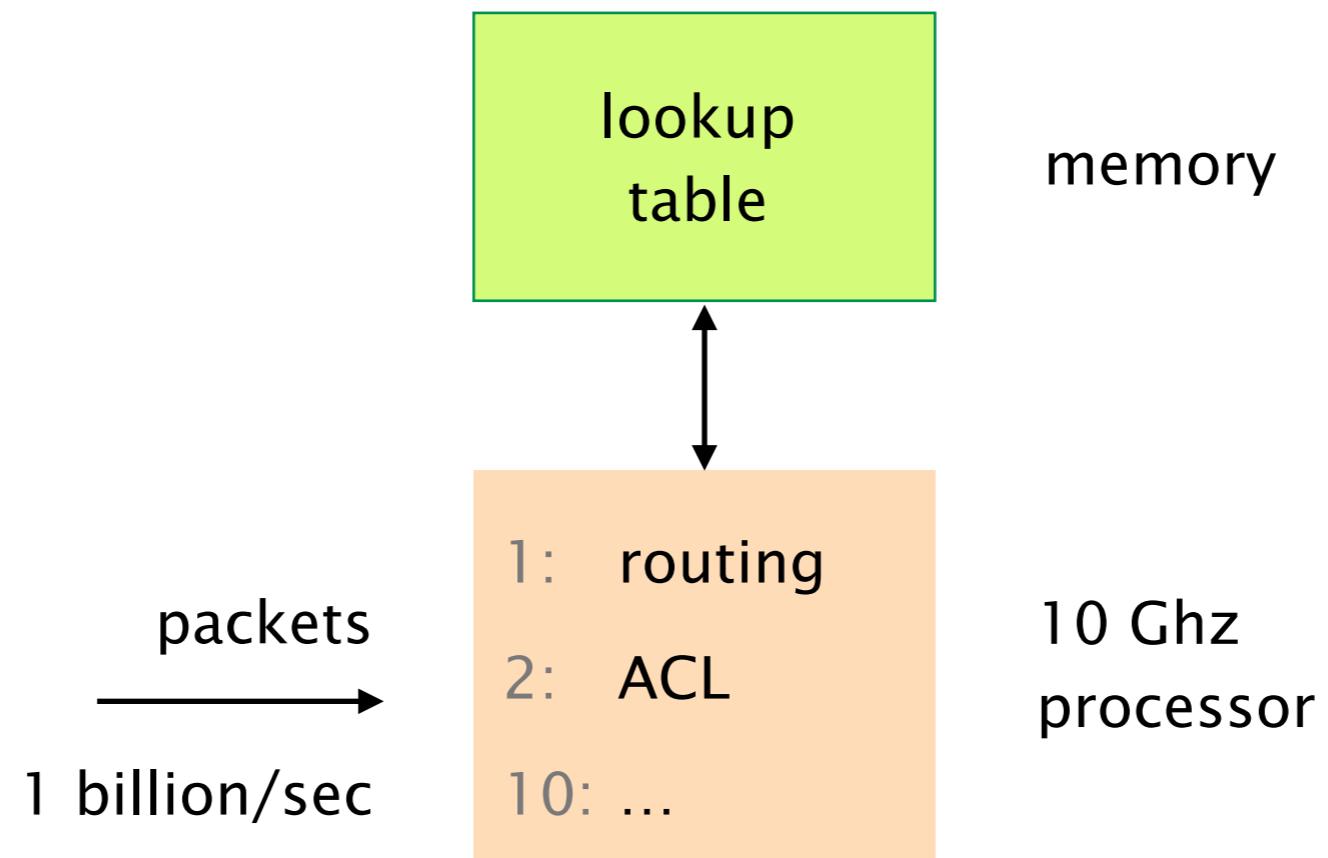
Outline

- Conventional switch chip are inflexible
- SDN demands flexibility...sounds expensive...
- **How do we do it: The RMT switch model**
- Flexibility costs less than 15%

What kind of switch architecture could support flexibility and yet run at Terabits per second?

Throughput aggregate	1 Tbps
Packet size average	1000 bits
# operations per packet (avg.)	10
Requirements	10 billion op./second

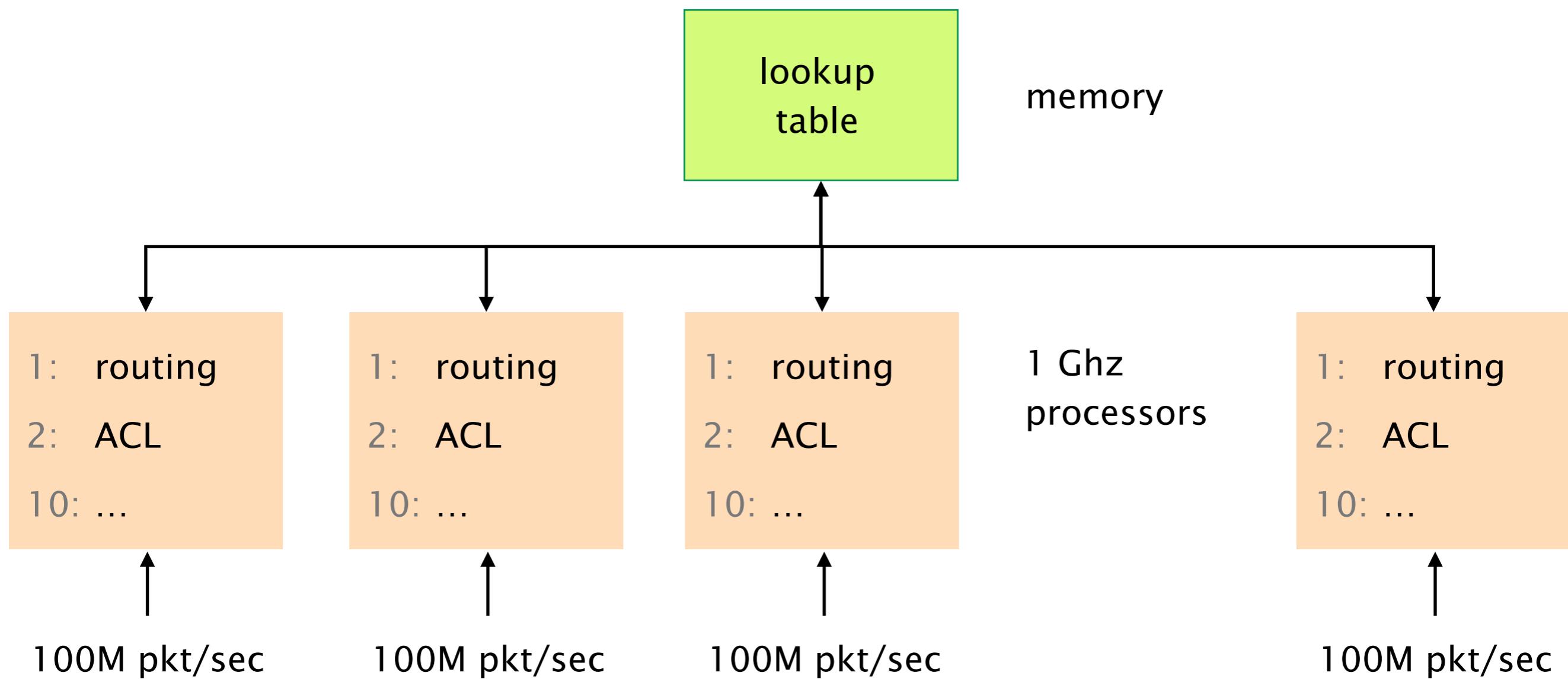
If our switch has a single processor,
this would require us to run it at **10 Ghz**...
not feasible



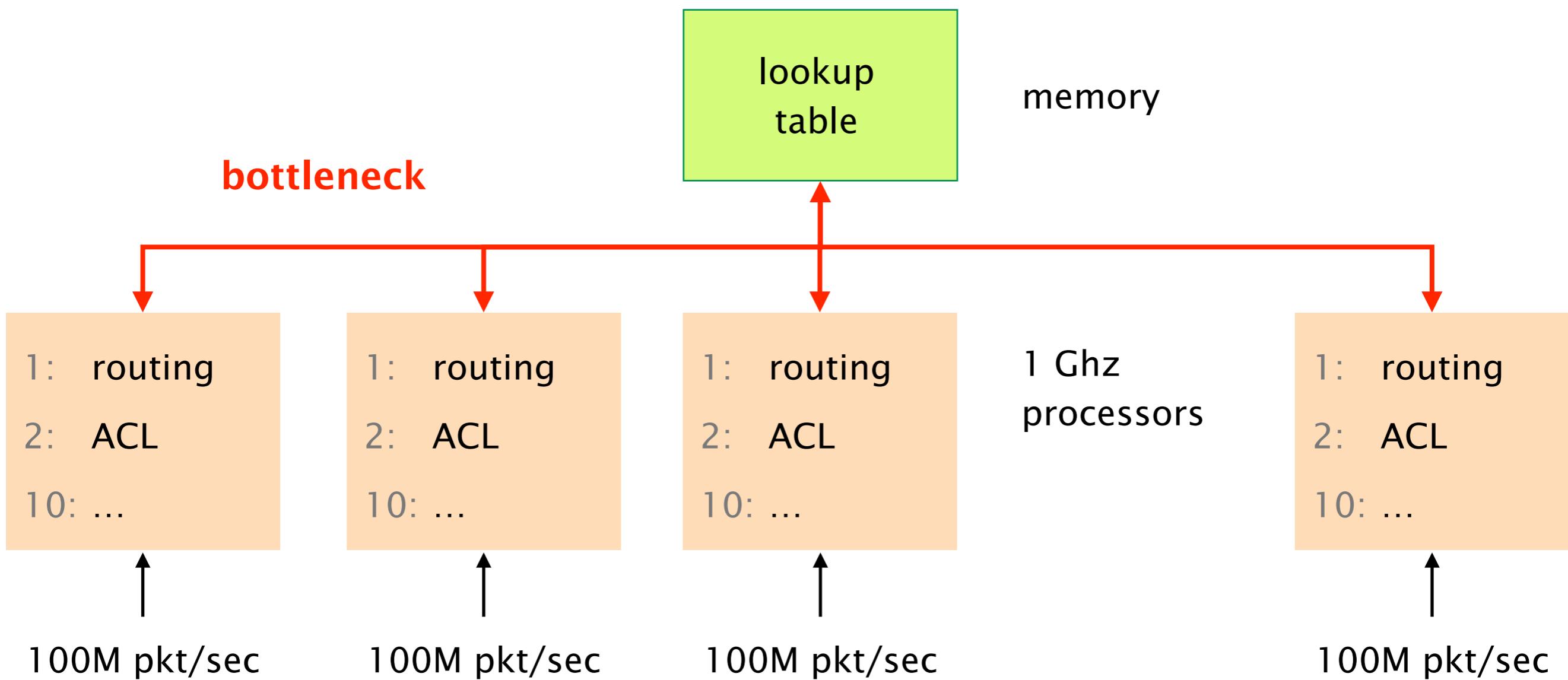
Let's parallelize things with a
packet-parallel architecture

What about we duplicate the processing units?

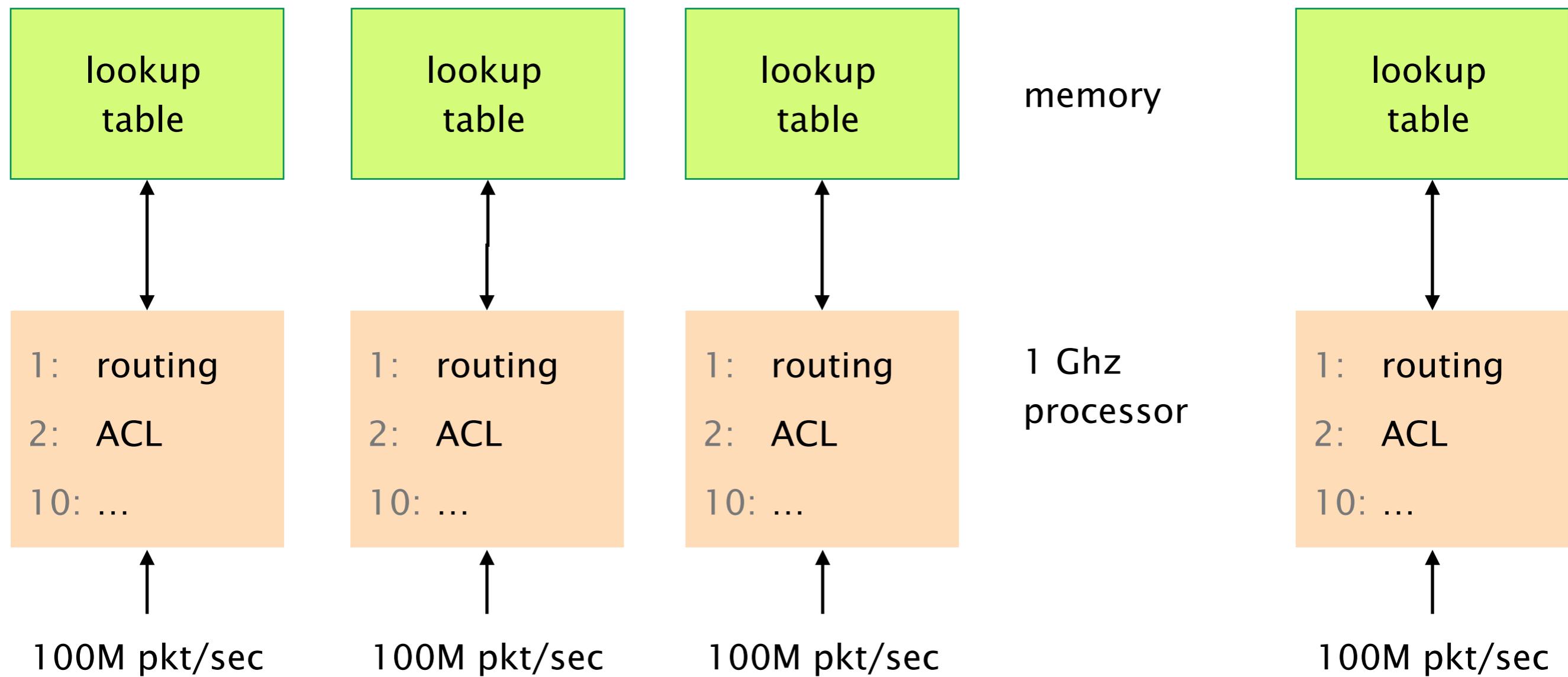
Each of which clocked at 1 Ghz



The next issue is to scale the
the memory-to-CPU bandwidth

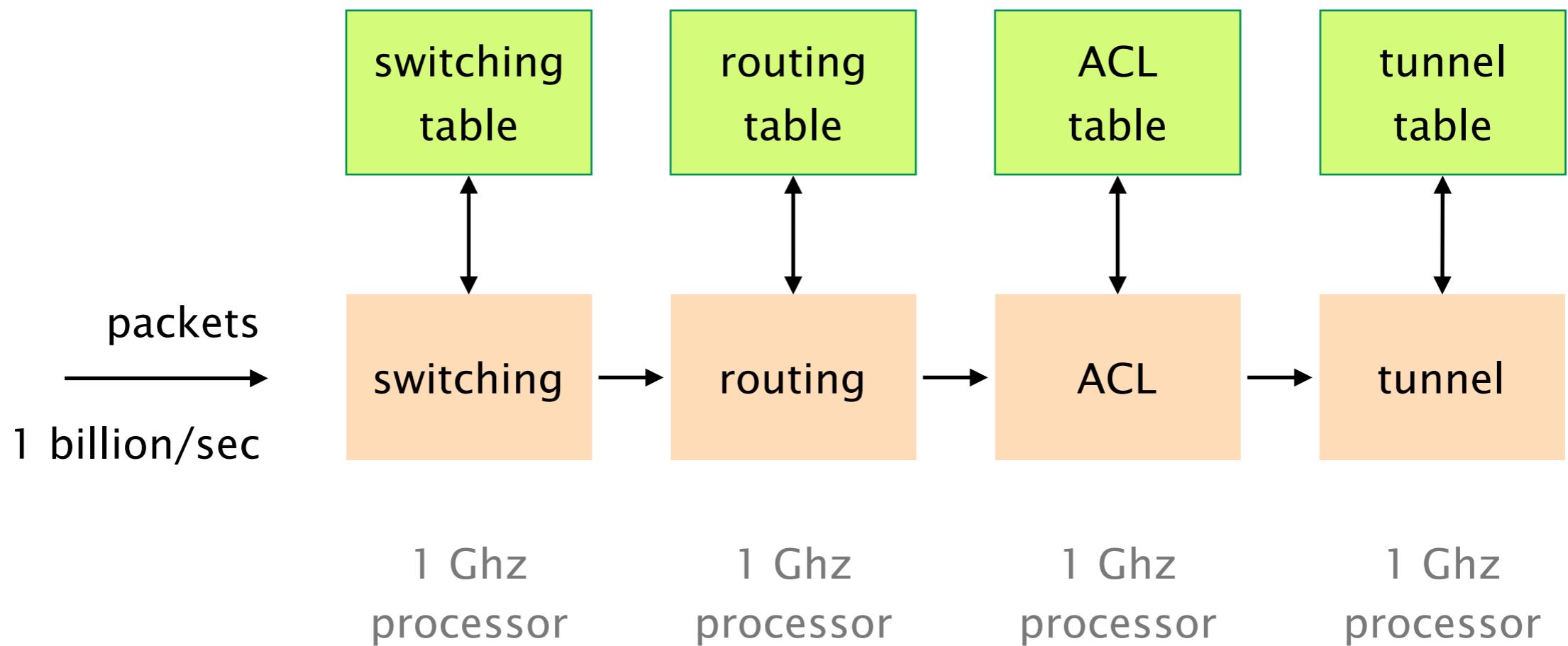


We could replicate the memory of course...
but that comes at **a huge cost** in die area

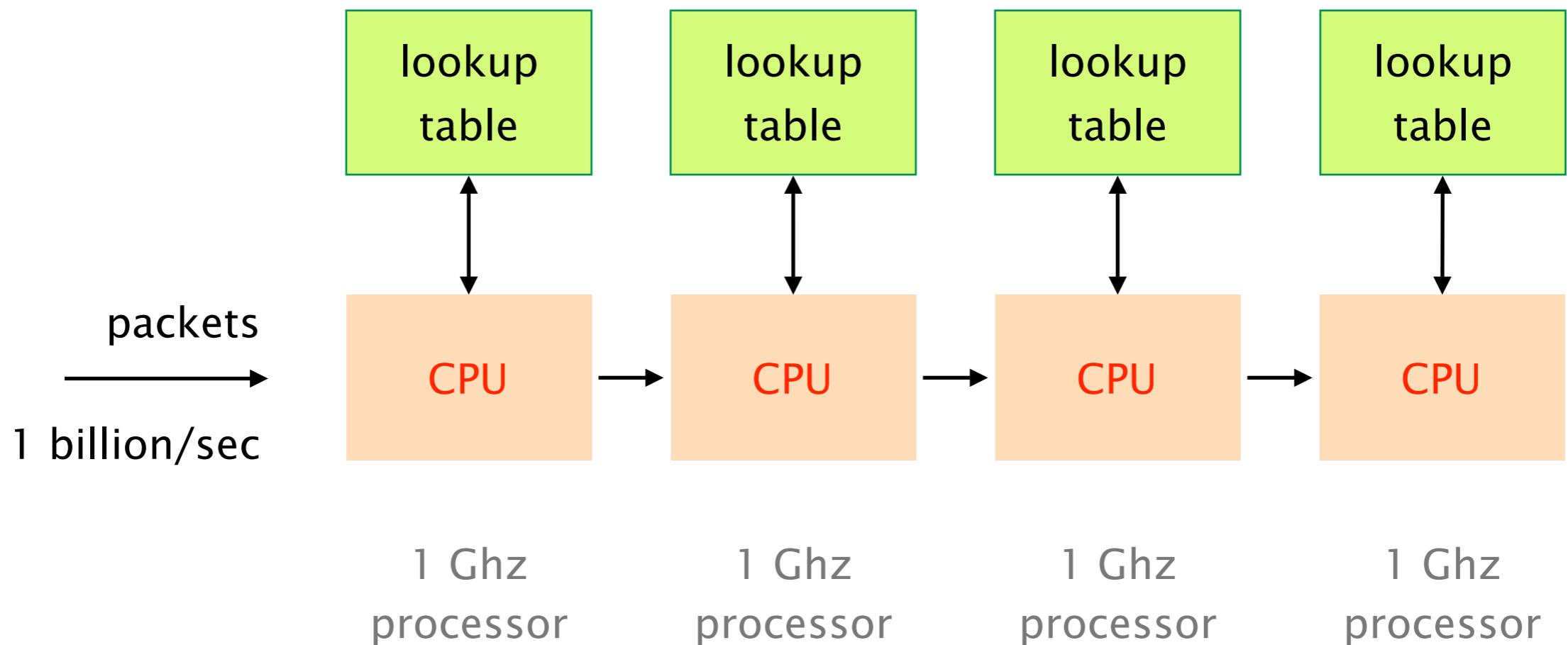


What if we organize the processing
as a **pipeline** instead?

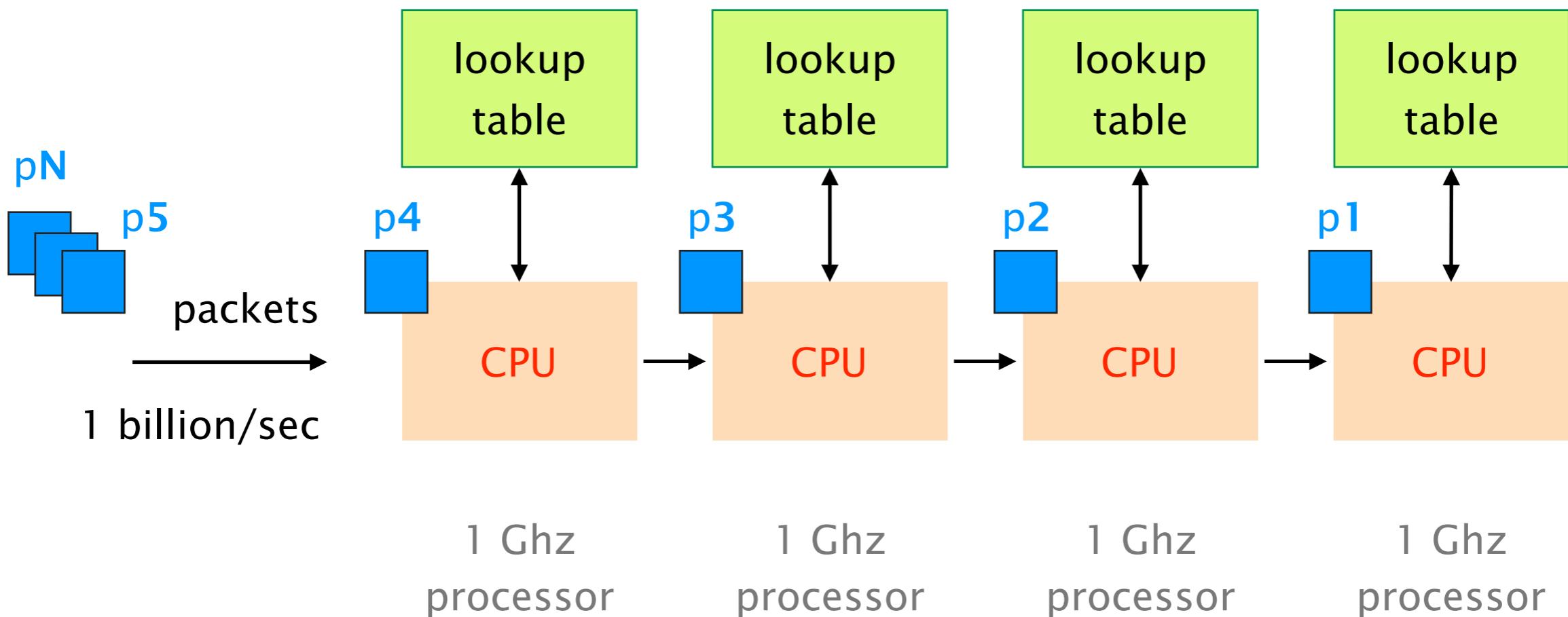
Pipelined architectures organize processing through a sequence of processing units and local memory



For flexibility,
each processing unit/memory can be made generic



Each CPU can process distinct packets, with up to 10 packets going through the pipeline simultaneously



The runtime behavior of the parser & the match stages is defined through the RMT abstract model

The RMT Abstract Model

- Parse graph
- Table graph

How do we implement in hardware
a programmable parser and a logical pipeline?

How do we implement in hardware a programmable parser and a logical pipeline?

The screenshot shows the first page of a PDF document titled 'Design Principles for Packet Parsers'. The page includes the authors' names (Glen Gibb, George Varghese, Mark Horowitz, Nick McKeown), their institutions (Stanford University, Microsoft Research), and their email addresses. The abstract discusses the challenges of parsing packet headers in network devices. The categories and subject descriptors mention Computer-Communication Networks, Network Architecture and Design, and Network Communications. The keywords include Parsing, Design principles, and Reconfigurable parsers. The introduction section is partially visible at the bottom.

Design Principles for Packet Parsers

Glen Gibb[†], George Varghese[‡], Mark Horowitz[†], Nick McKeown[†]
†Stanford University ‡Microsoft Research
`{grg, horowitz, nickm}@stanford.edu` `varghese@microsoft.com`

ABSTRACT

All network devices must parse packet headers to decide how packets should be processed. A 64×10 Gb/s Ethernet switch must parse one billion packets per second to extract fields used in forwarding decisions. Although a necessary part of all switch hardware, very little has been written on parser design and the trade-offs between different designs. Is it better to design one fast parser, or several slow parsers? What is the cost of making the parser reconfigurable in the field? What design decisions most impact power and area?

In this paper, we describe trade-offs in parser design, identify design principles for switch and router designers, and describe a parser generator that outputs synthesizable Verilog that is available for download. We show that i) packet parsers today occupy about 1-2% of the chip, and ii) while future packet parsers will need to be programmable, this only doubles the (already small) area needed.

Categories and Subject Descriptors

C.2.1 [Computer-Communication Networks]: Network Architecture and Design—*Network Communications*

Keywords

Parsing; Design principles; Reconfigurable parsers

1. INTRODUCTION

Despite their variety, *every* network device examines fields

The diagram illustrates the structure of a TCP packet. It starts with an 'Ethernet' layer (Len: 14B) followed by an 'IPv4' layer (Len: ?). The 'IPv4' layer contains a 'Next: IPv4' header and a payload of 20B. This is followed by a 'TCP' layer (Len: ?) containing a 'Next: TCP' header and a payload of 20B. Finally, there is a 'Payload' layer. Arrows indicate the flow from left to right through the layers.

Figure 1: A TCP packet.

In practice, packets often contain many more headers. These extra headers carry information about higher level protocols (e.g., HTTP headers) or additional information that existing headers do not provide (e.g., VLANs¹ in a college campus, or MPLS² in a public Internet backbone). It is common for a packet to have eight or more different packet headers during its lifetime.

To parse a packet, a network device has to identify the headers in sequence before extracting and processing specific fields. A packet parser seems straightforward since it knows *a priori* which header types to expect.

In practice, designing a parser is quite challenging:

1. **Throughput.** Most parsers must run at line-rate, supporting continuous minimum-length back-to-back packets. A 10 Gb/s Ethernet link can deliver a new packet every 70 ns; a state-of-the-art Ethernet switch ASIC with 64×40 Gb/s ports must process a new packet every 270 ps.

[ANCS'13]

Parsing is the (complex) process of identifying and extracting the appropriate fields in a packet header

Throughput

Parser must run at line-rate
parse 1 packet every 70 ns on a 10 Gbps link

Dependency

Parsing involves sequential processing
as headers typically point to the next one

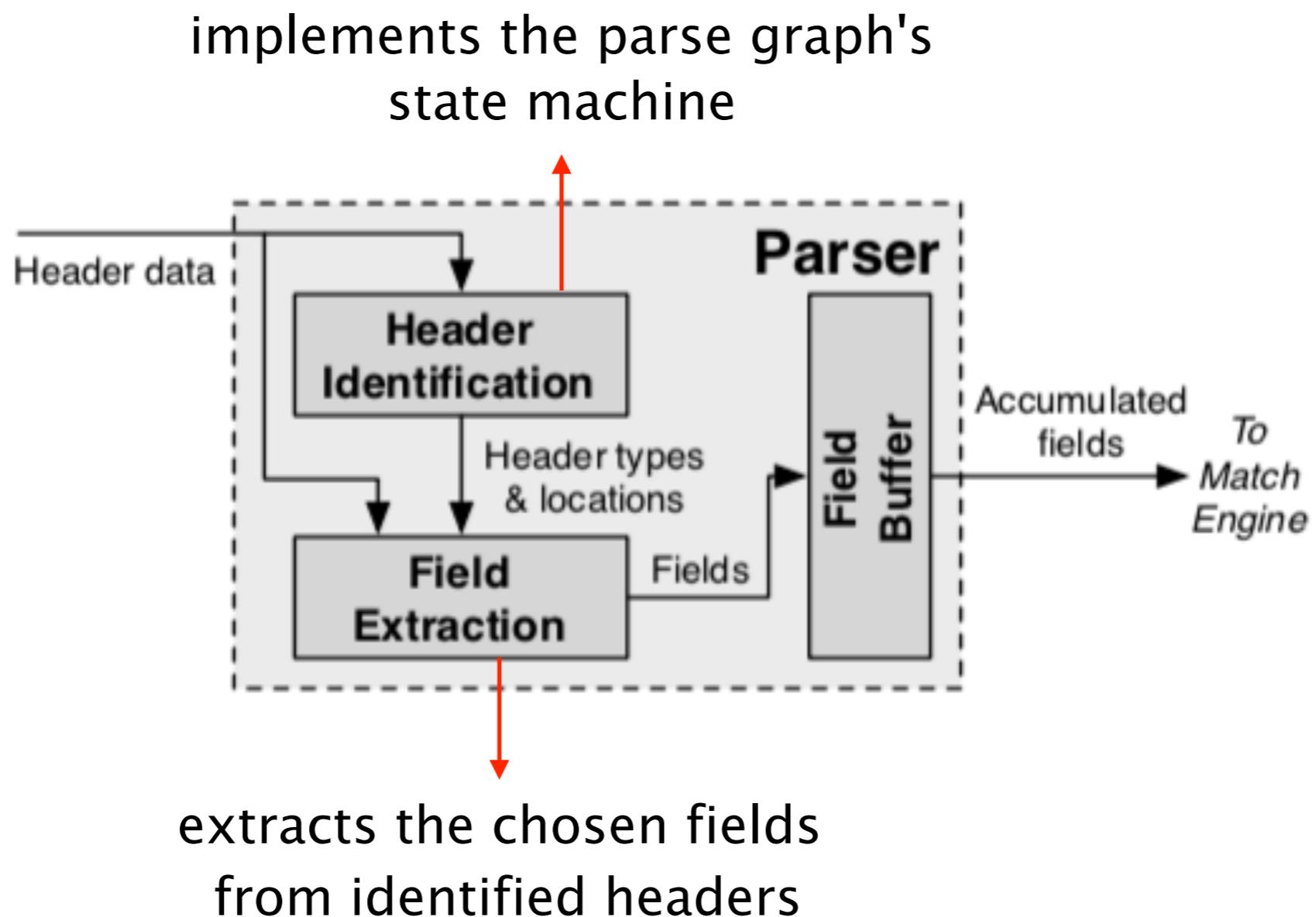
Incompleteness

Some headers do not even identify
the subsequent header

Heterogeneity

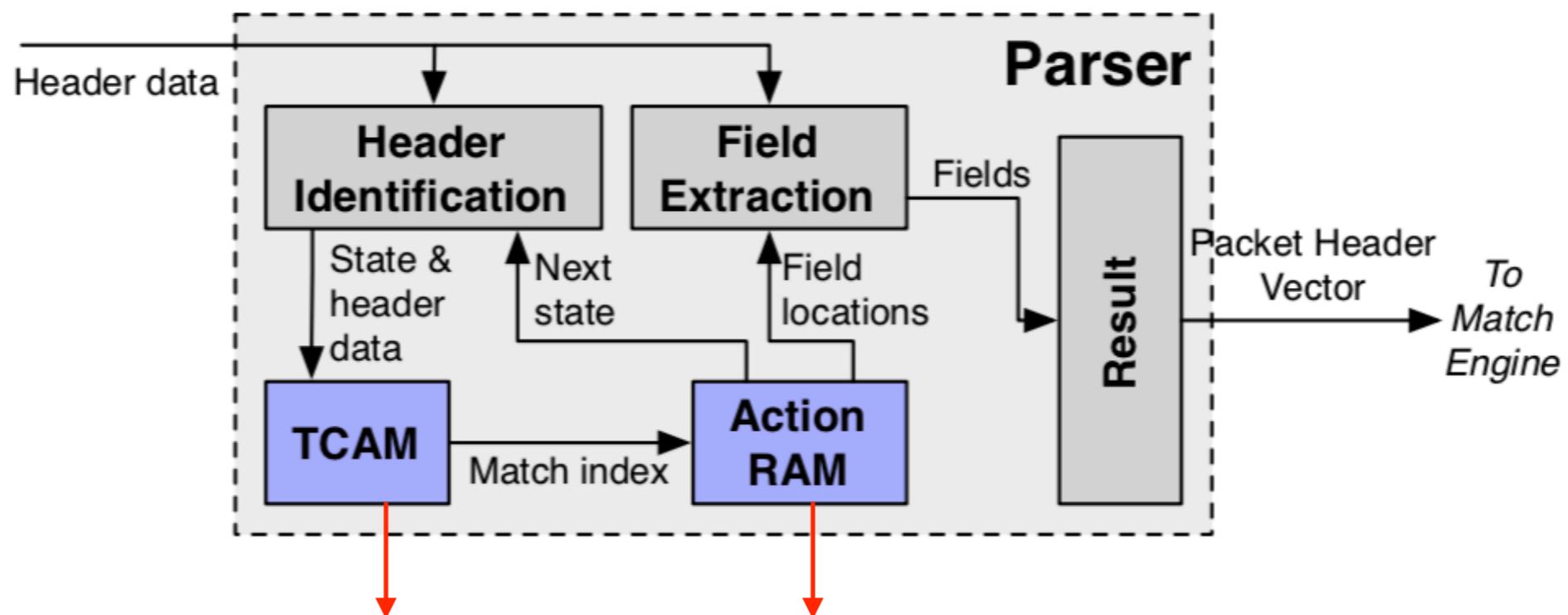
Many header formats exist that
can appear in various orders/locations

A parser can be divided into two separate blocks:
header identification and field extraction



In a programmable parser, the two modules rely on
runtime information instead of hard-coded logic

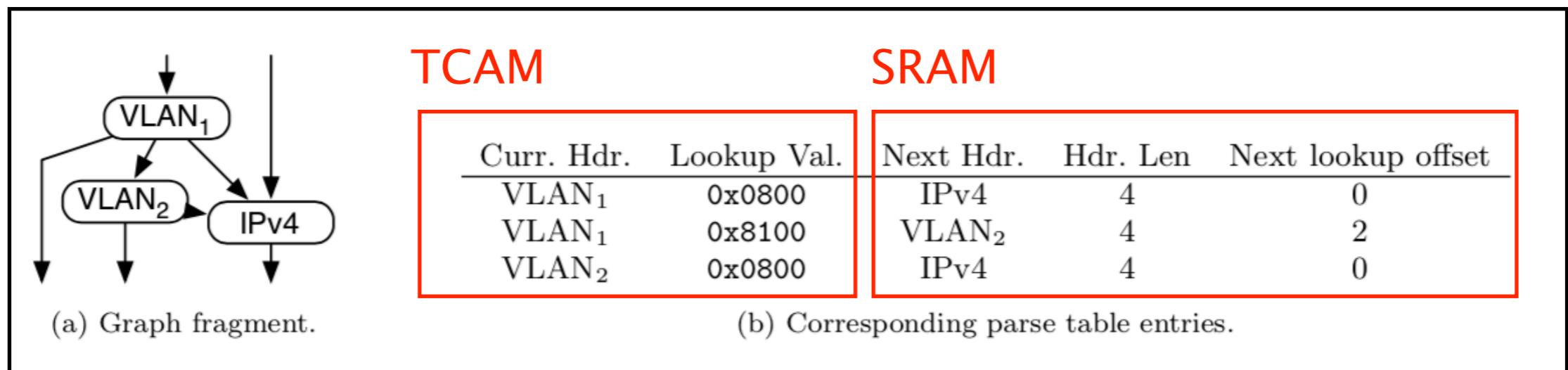
stored in memory,
e.g. in SRAM and/or TCAM



stores the bit sequences
that identify the headers

stores the next state,
the fields to extract,
and any other data (if any)

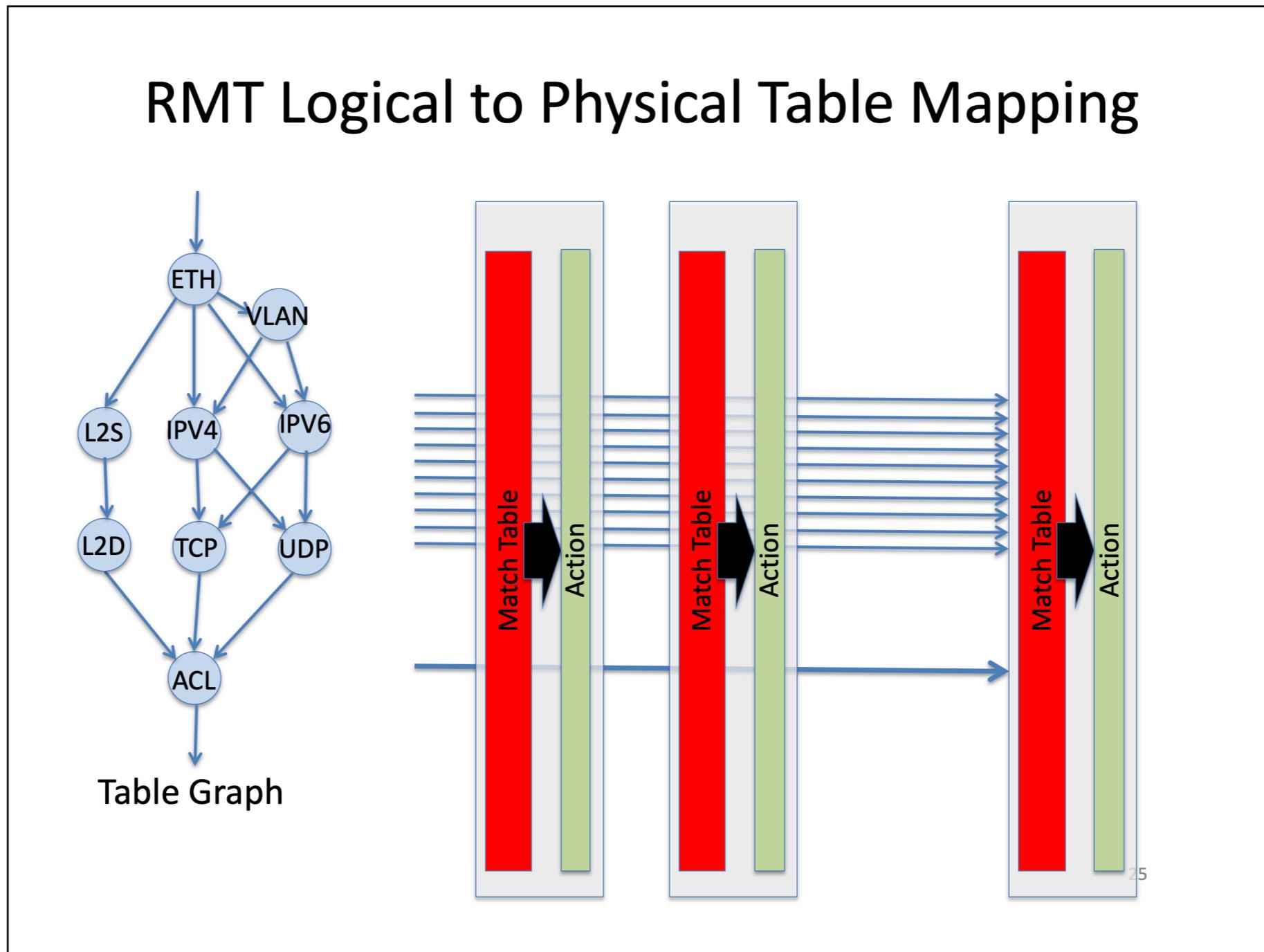
Linked together, a SRAM and TCAM can encode the transition table attached to a parsing graph



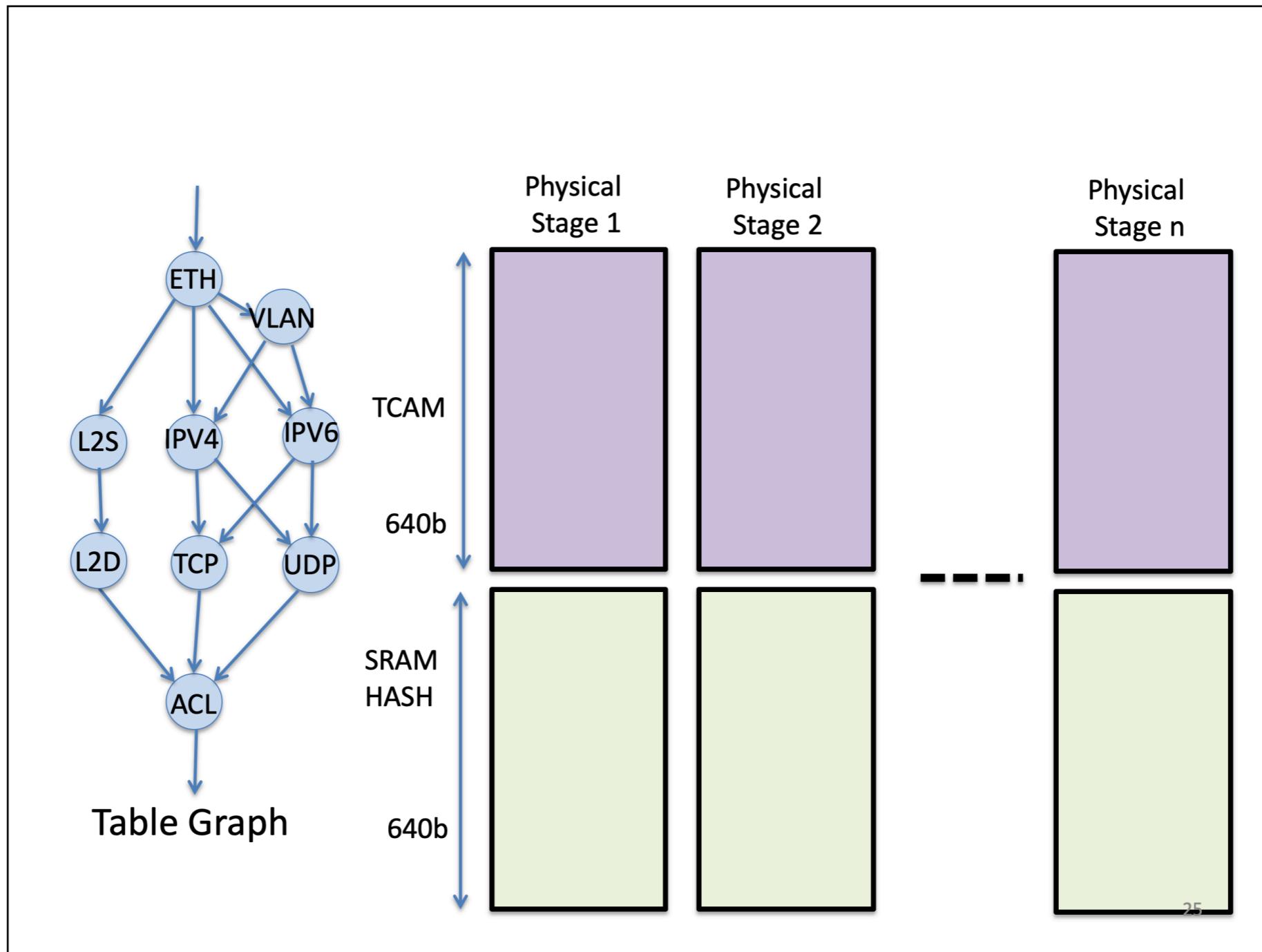
Source: Design Principles for Packet Parsers, Gibb et al.

How do we implement in hardware
a programmable parser and a logical pipeline?

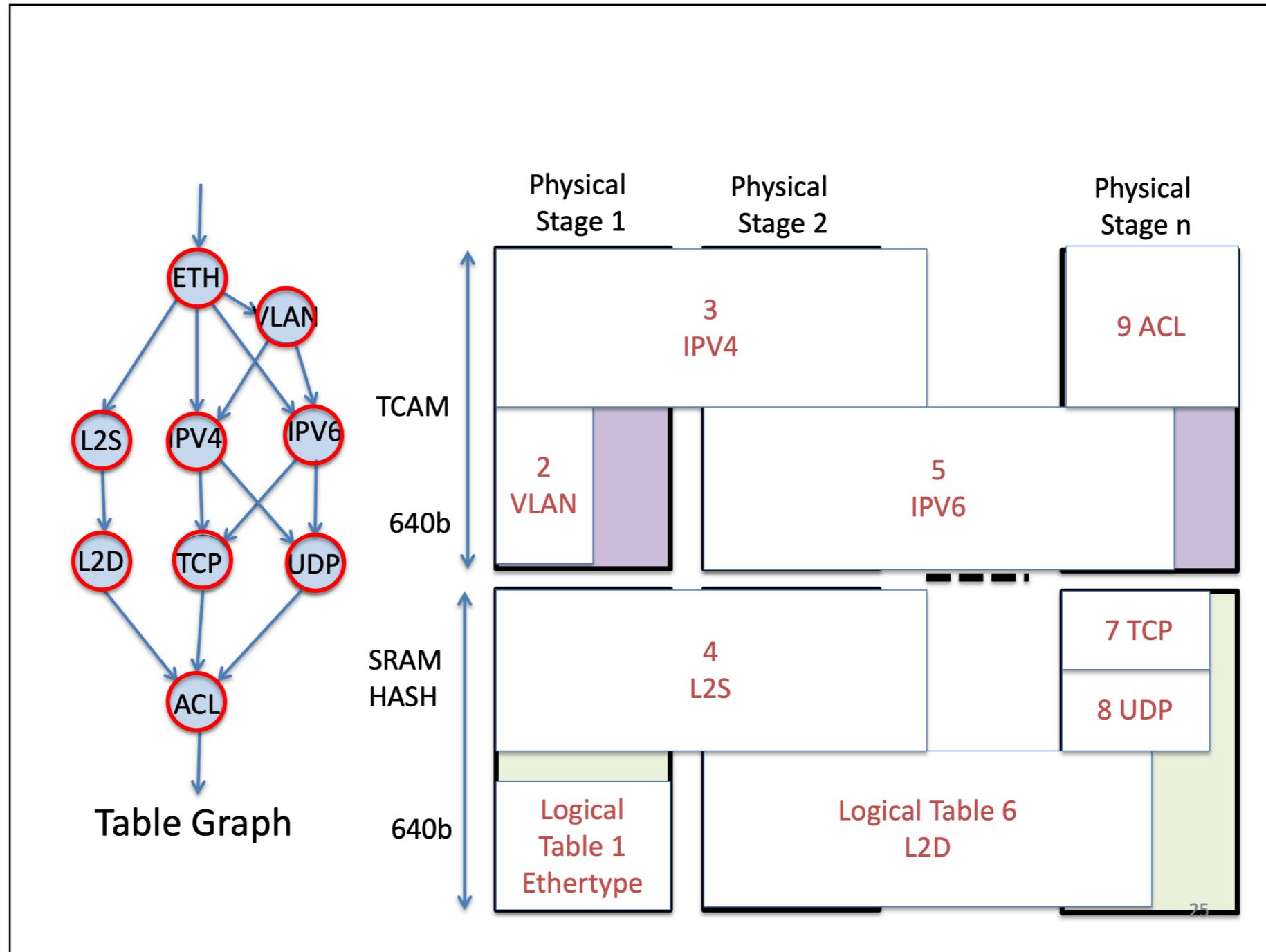
A compiler translates a given RMT logical pipeline (specified in P4) into a physical one



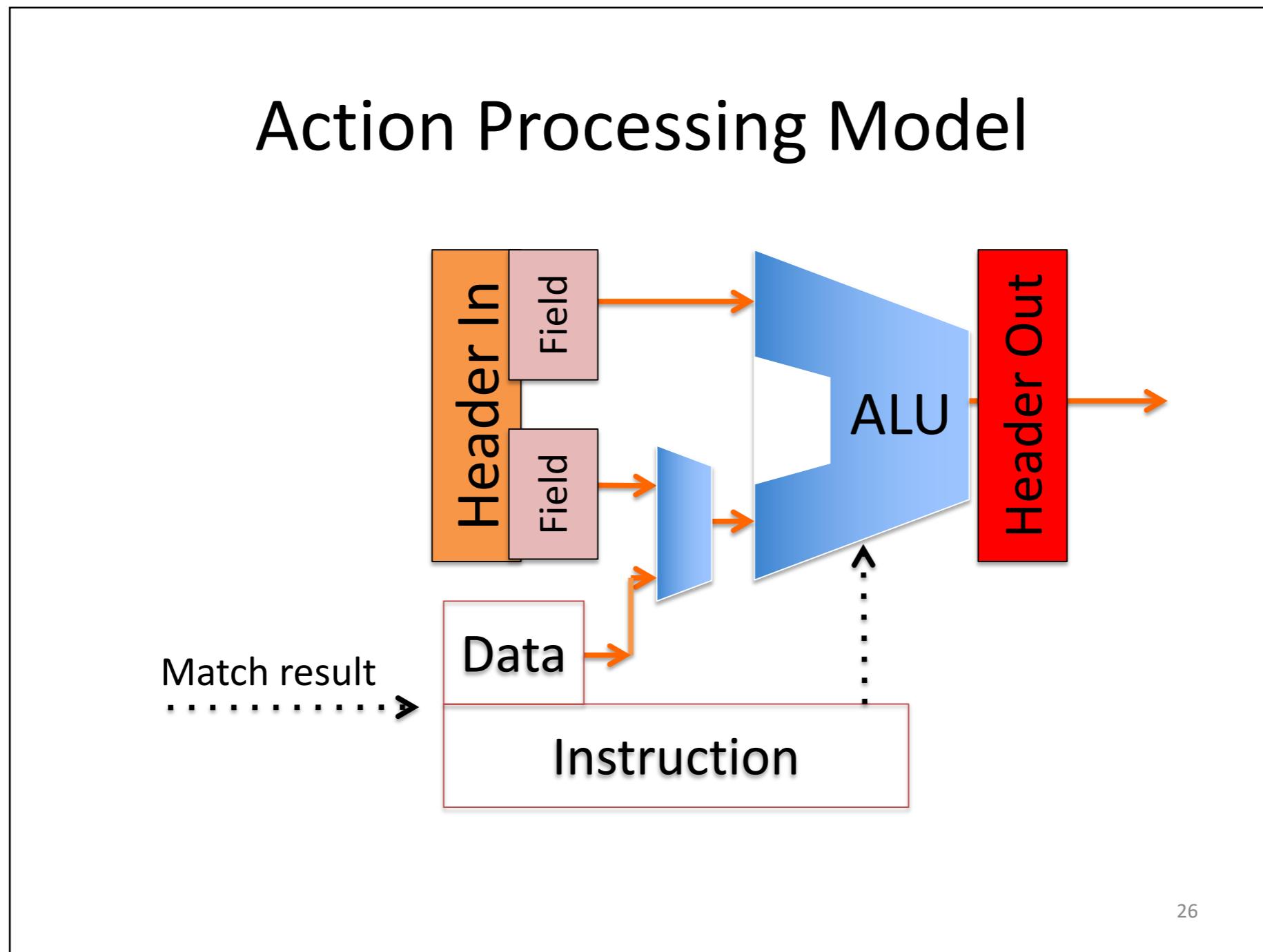
Each physical stage contains dedicated SRAM,
for exact matches, and TCAM, for ternary matches



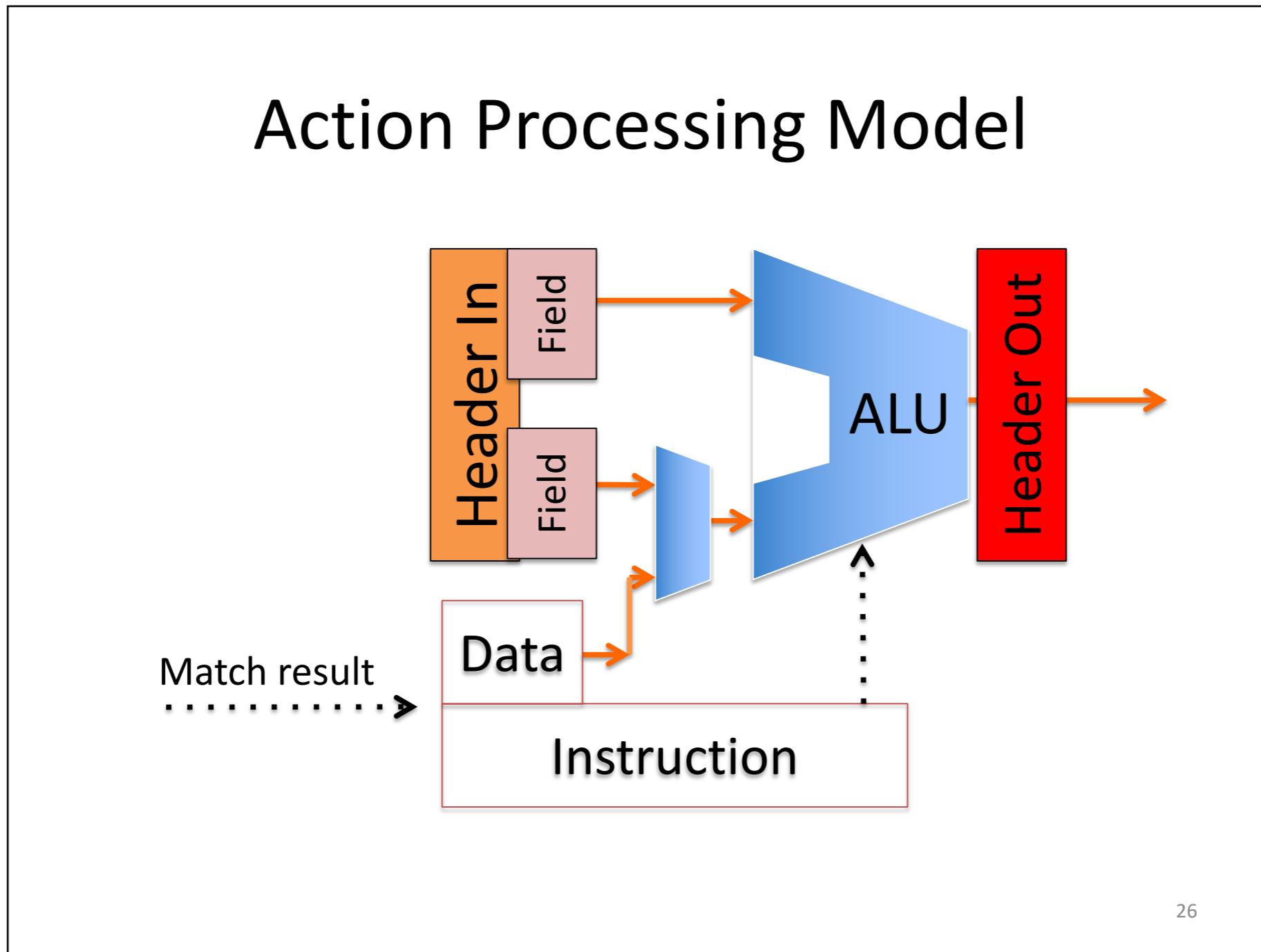
The compiler maps each individual logical stage to one or more physical stage.



The RMT pipeline relies on many Arithmetic Logic Units (ALU) to perform actions on the result of a match

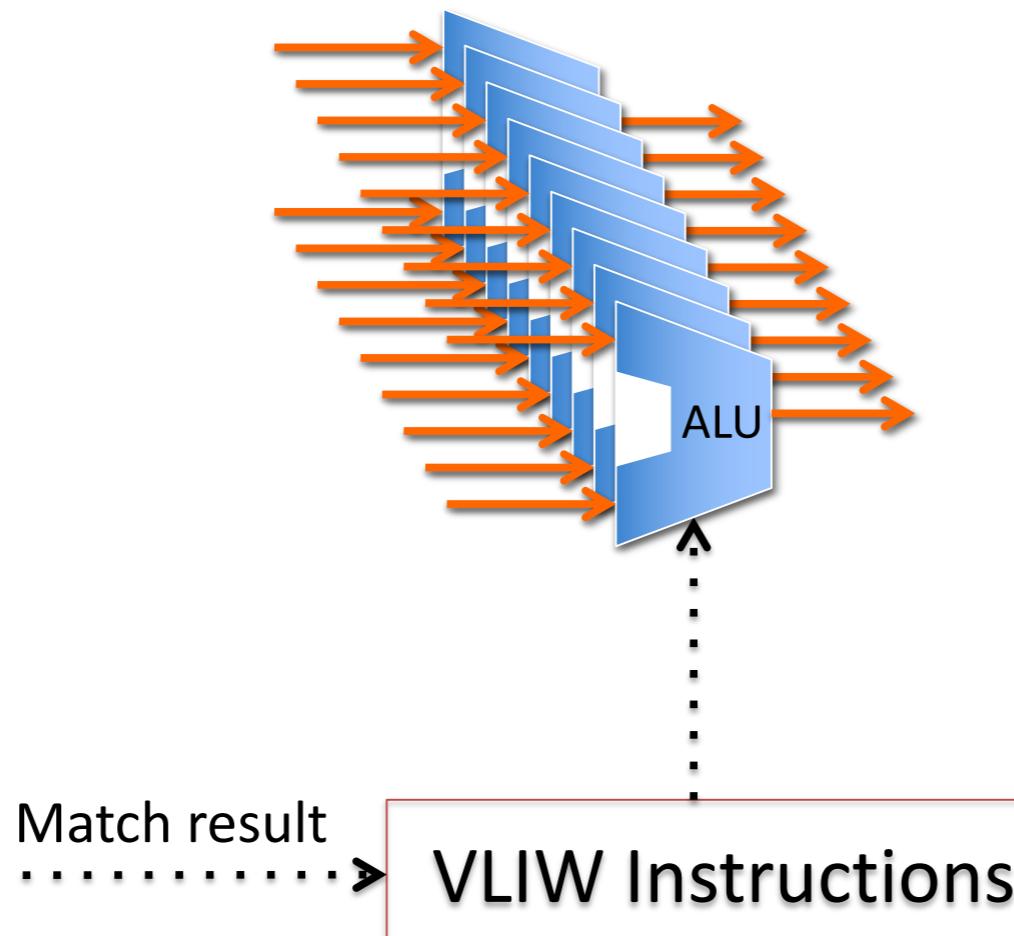


Each ALU modifies only one word of a header
(a header is composed of *many* words)



Each stage of the RMT pipeline contains one ALU per word of the header vector (that's *a lot* of ALUs)

Modeled as Multiple VLIW CPUs per Stage



The RMT pipeline in a few statistics

Our Switch Design

- 64 x 10Gb ports
 - 960M packets/second
 - 1GHz pipeline
- Programmable parser
- 32 Match/action stages
- Huge TCAM: 10x current chips
 - 64K TCAM words x 640b
- SRAM hash tables for exact matches
 - 128K words x 640b
- 224 action processors per stage
- All OpenFlow statistics counters

Building a RMT pipeline is **only 15% more expensive**
than building a fixed-function switching pipeline

Outline

- Conventional switch chip are inflexible
- SDN demands flexibility...sounds expensive...
- How do I do it: The RMT switch model
- **Flexibility costs less than 15%**

In terms of die area, flexibility is not very expensive
at least, not anymore... mainly thanks to Moore's law

Chip Comparison with Fixed Function Switches

Area

Section	Area % of chip	Extra Cost
IO, buffer, queue, CPU, etc	37%	0.0%
Match memory & logic	54.3%	8.0%
VLIW action engine	7.4%	5.5%
Parser + deparser	1.3%	0.7%
Total extra area cost		14.2%



The same lesson applies for power

Chip Comparison with Fixed Function Switches

Area		
Section	Area % of chip	Extra Cost
IO, buffer, queue, CPU, etc	37%	0.0%
Match memory & logic	54.3%	8.0%
VLIW action engine	7.4%	5.5%
Parser + deparser	1.3%	0.7%
Total extra area cost		14.2%

Power		
Section	Power % of chip	Extra Cost
I/O	26.0%	0.0%
Memory leakage	43.7%	4.0%
Logic leakage	7.3%	2.5%
RAM active	2.7%	0.4%
TCAM active	3.5%	0.0%
Logic active	16.8%	5.5%
Total extra power cost		12.4%

Conclusion

- How do we design a flexible chip?
 - The RMT switch model
 - Bring processing close to the memories:
 - pipeline of many stages
 - Bring the processing to the wires:
 - 224 action CPUs per stage
- How much does it cost?
 - 15%
- Lots of the details how we designed this in 28nm CMOS are in the paper

P4 hardware
target

Probabilistic
data structures

Can we trade some accuracy
for memory?



← You are looking at a stream of data.

In networking, we usually talk about **streams of packets**,
but these questions apply to other domains as well,
e.g. **search engines and databases**.



← You are looking at a stream of data (packets).

There are many questions you might ask:

Is a certain element (e.g. ip address) in the stream?

How frequently does an element appear?

How many distinct elements are in the stream?

What are the most frequent elements?



← You are looking at a stream of data (packets).

There are many questions you might ask:

Is a certain element (e.g. ip address) in the stream?

How frequently does an element appear?

How many distinct elements are in the stream?

What are the most frequent elements?

In P4, these questions are difficult to answer.



← You are looking at a stream of data (packets).
Today, I'll show you how set membership and frequency queries can be realized in P4.

PART 1

Is a certain element (e.g. ip address) in the stream?
→ Bloom filter

PART 2

How frequently does an element appear?
→ CountMin Sketch, Count Sketch, ...

part 1: set membership queries with Bloom filters

Is a certain element (e.g. ip address) in the stream?

(slides by Thomas Holterbach)

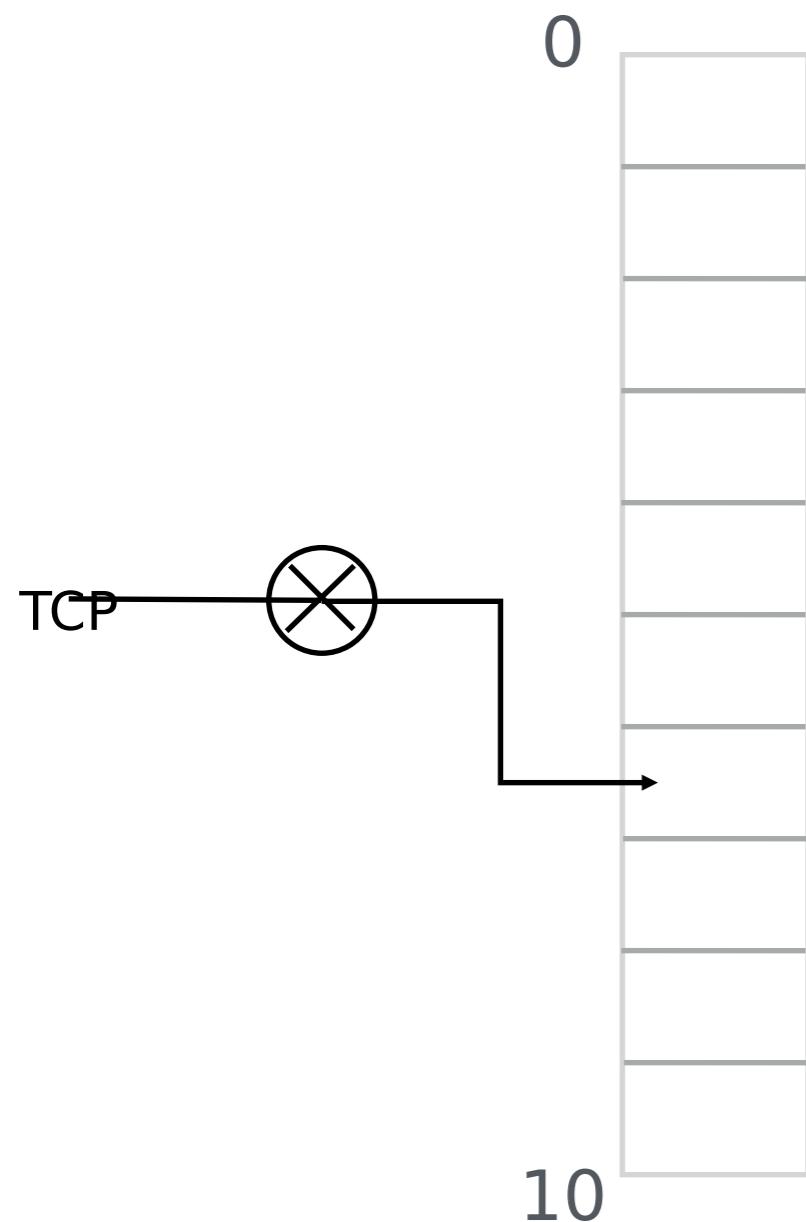
There are two common strategies
to implement a set

	strategy #1	strategy #2
output	Deterministic	
number of required operations	Probabilistic	

Intuitive implementation of a **set**

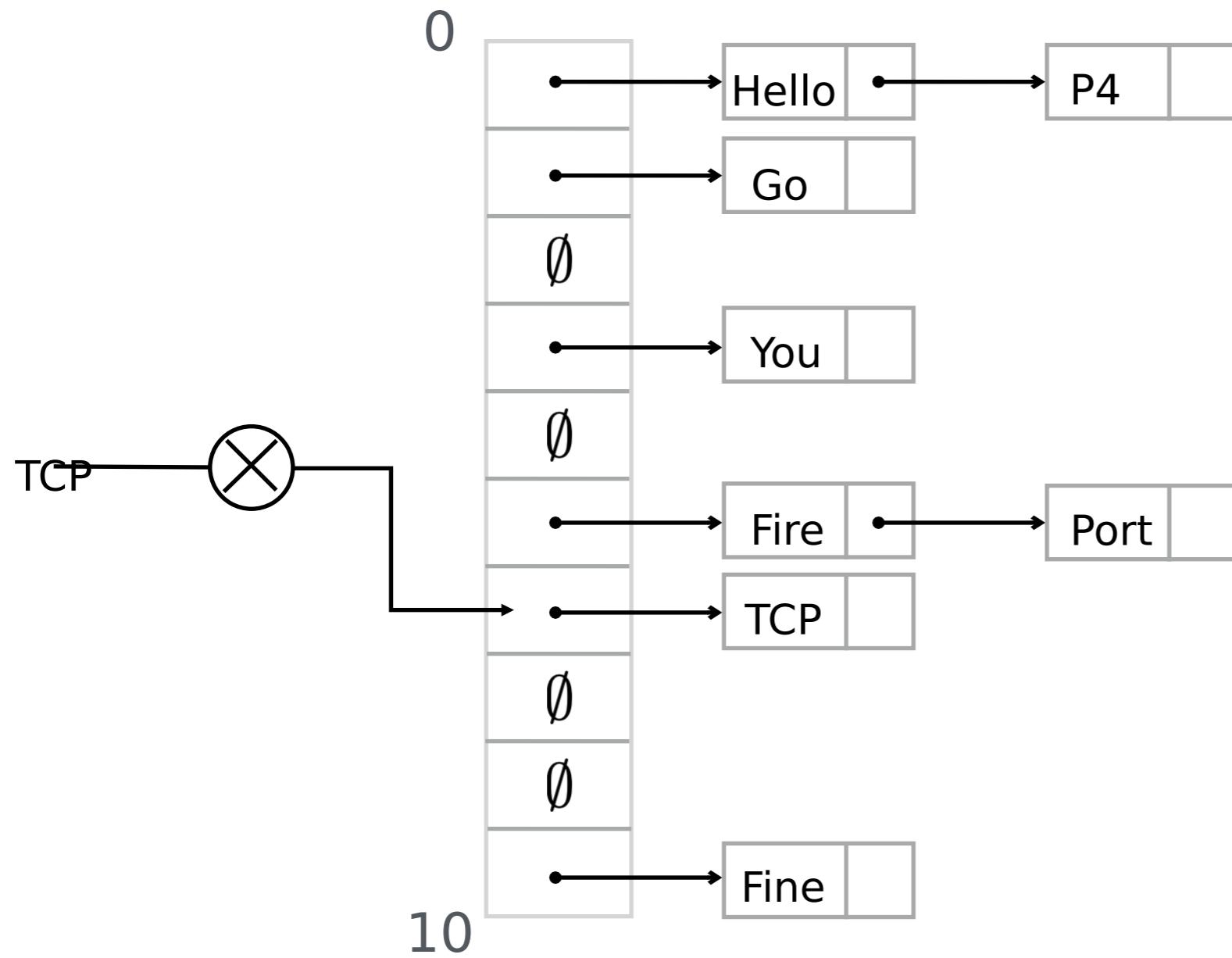
Intuitive implementation of a **set**

Separate-chaining



Intuitive implementation of a **set**

Separate-chaining



Intuitive implementation of a set

Separate-chaining

N elements and M cells

	list size
average	N/M
worse-case	N

Intuitive implementation of a **set**

Separate-chaining

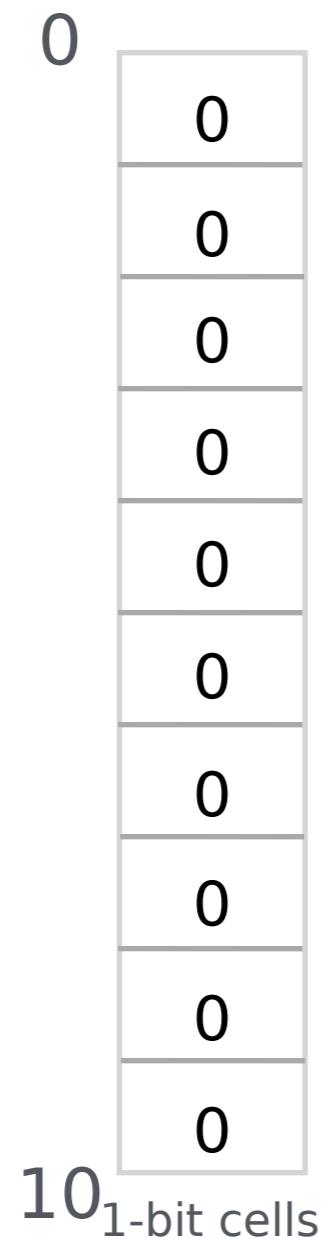
Pros: accurate and fast in the average case

Con: only works in hardware if there is a low number of elements (e.g. < 100)

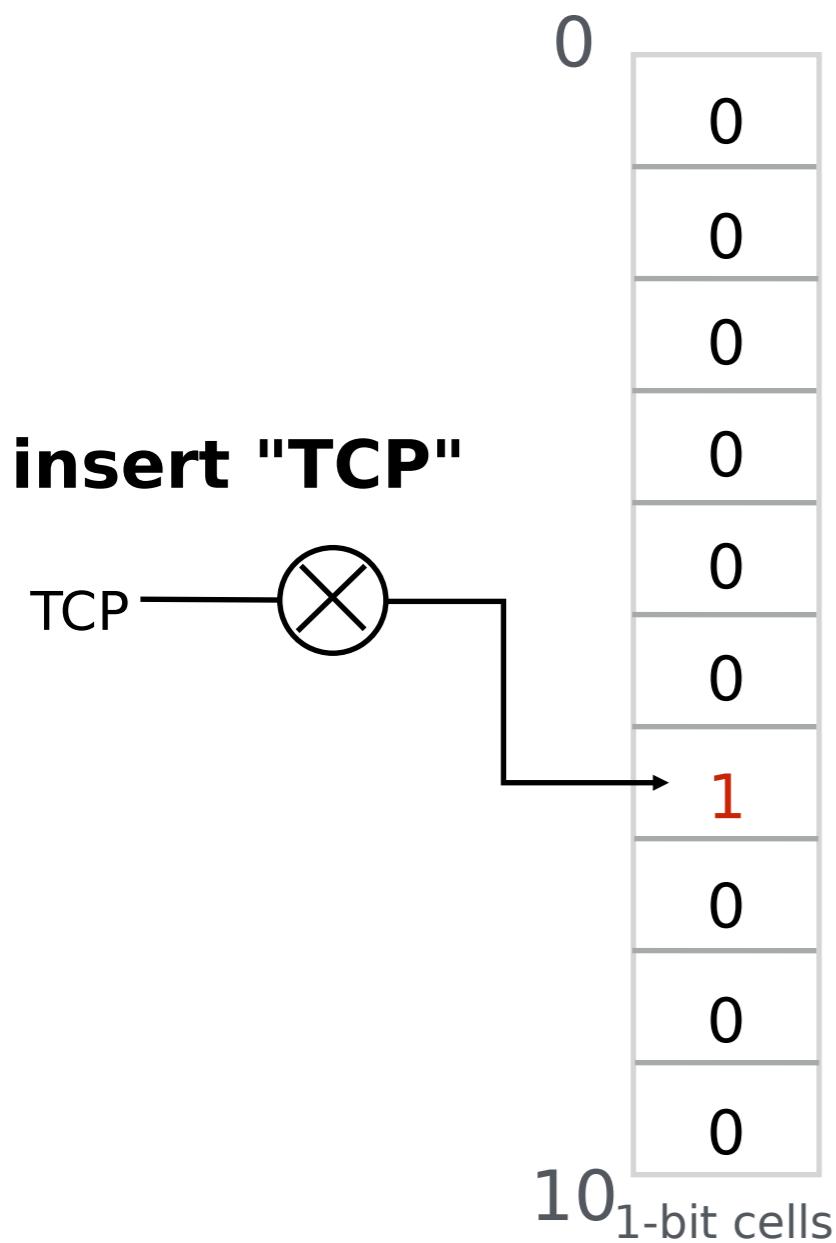
There are two common strategies
to implement a set

	strategy #1	strategy #2
output	Deterministic	Probabilistic
number of required operations	Probabilistic	Deterministic
<i>'probabilistic data structures'</i>		

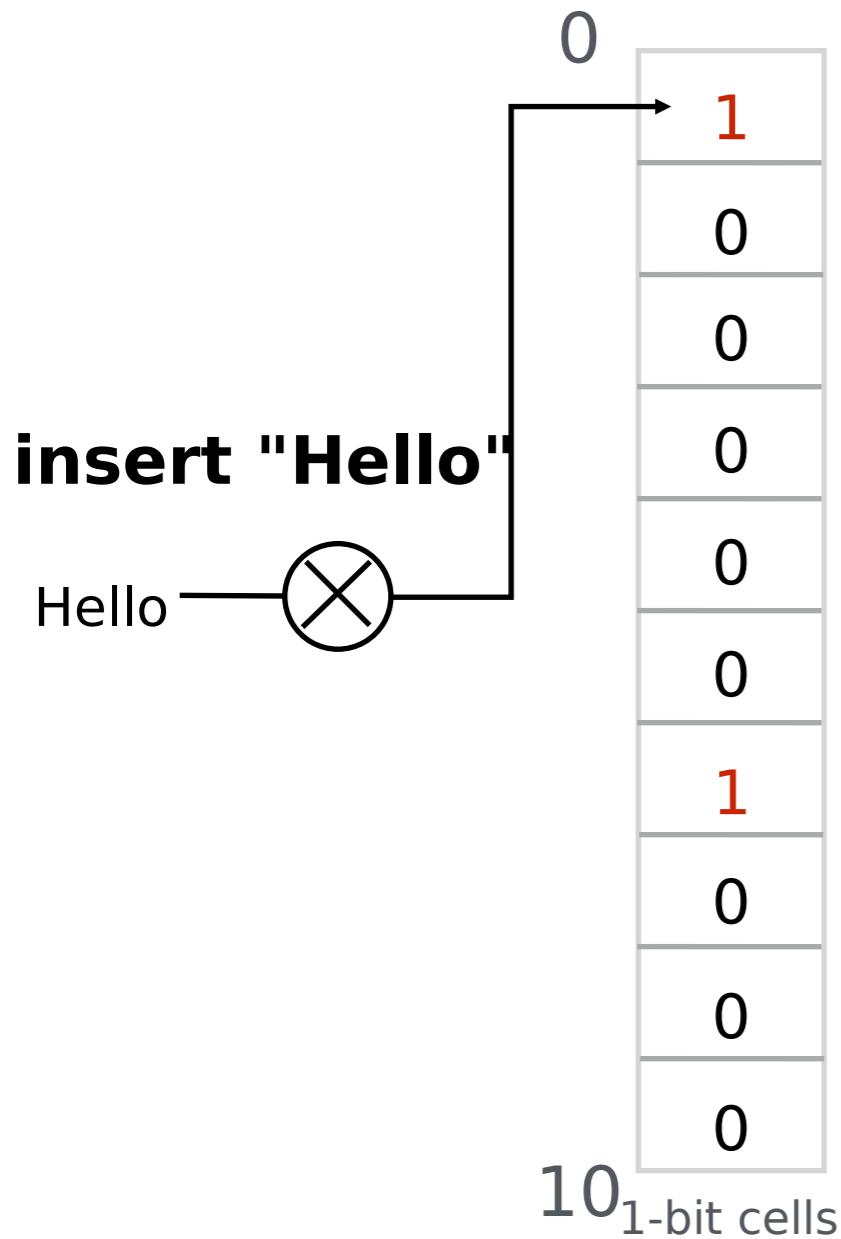
A simple approach for insertions and membership queries



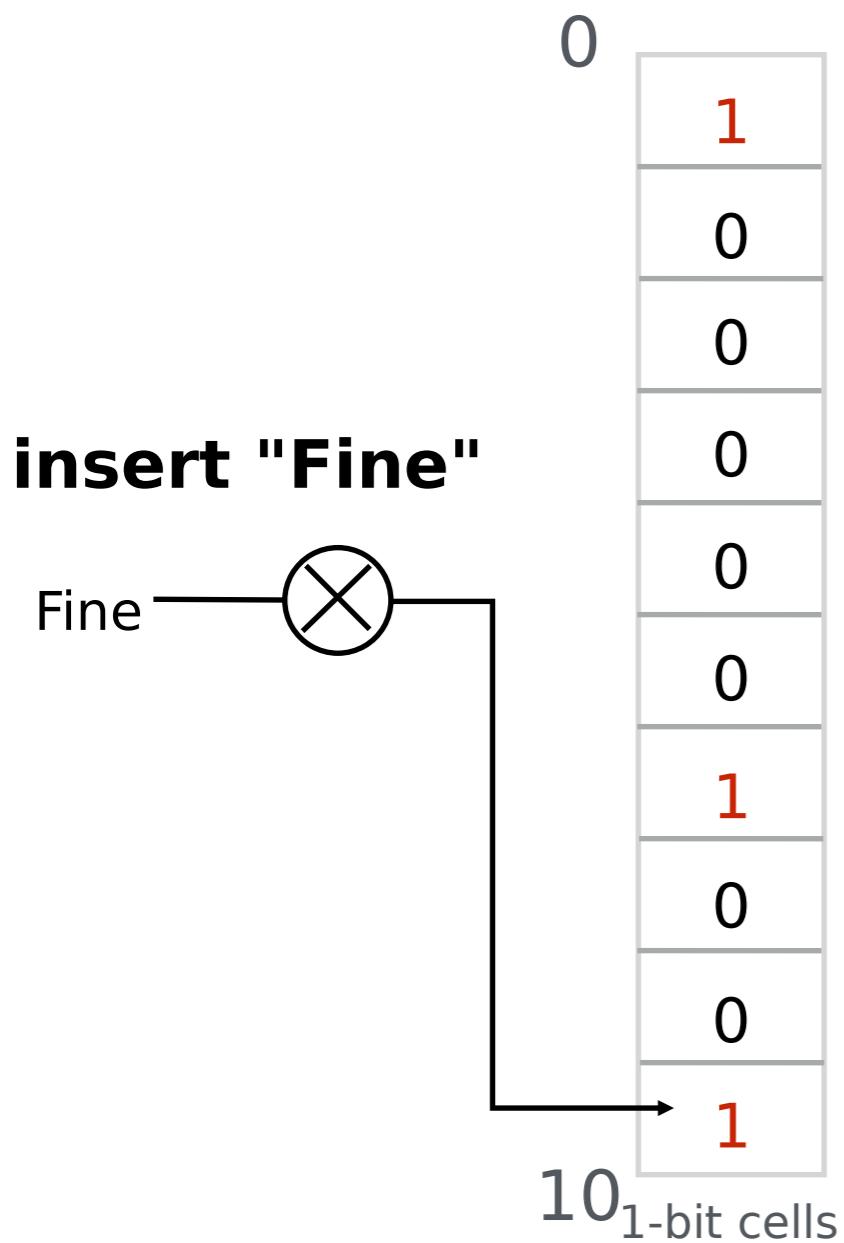
A simple approach for insertions and membership queries



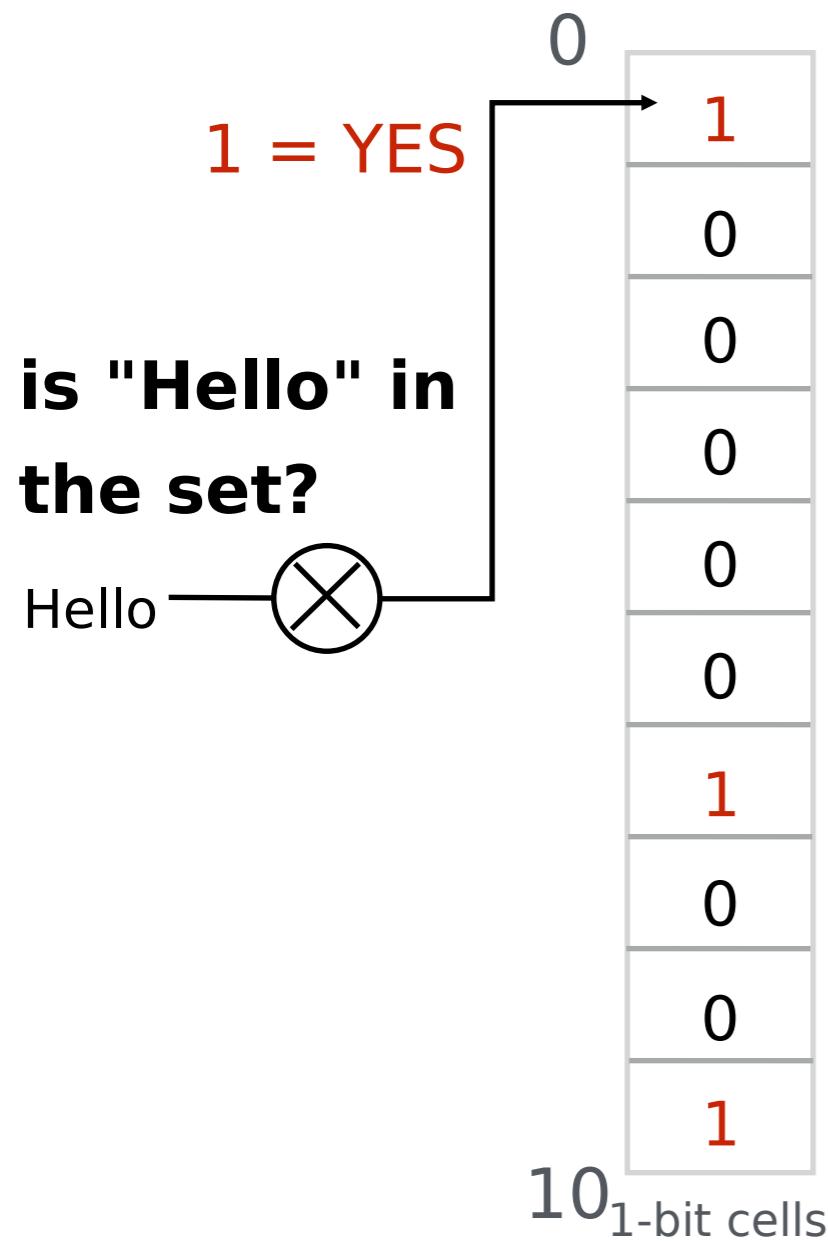
A simple approach for insertions and membership queries



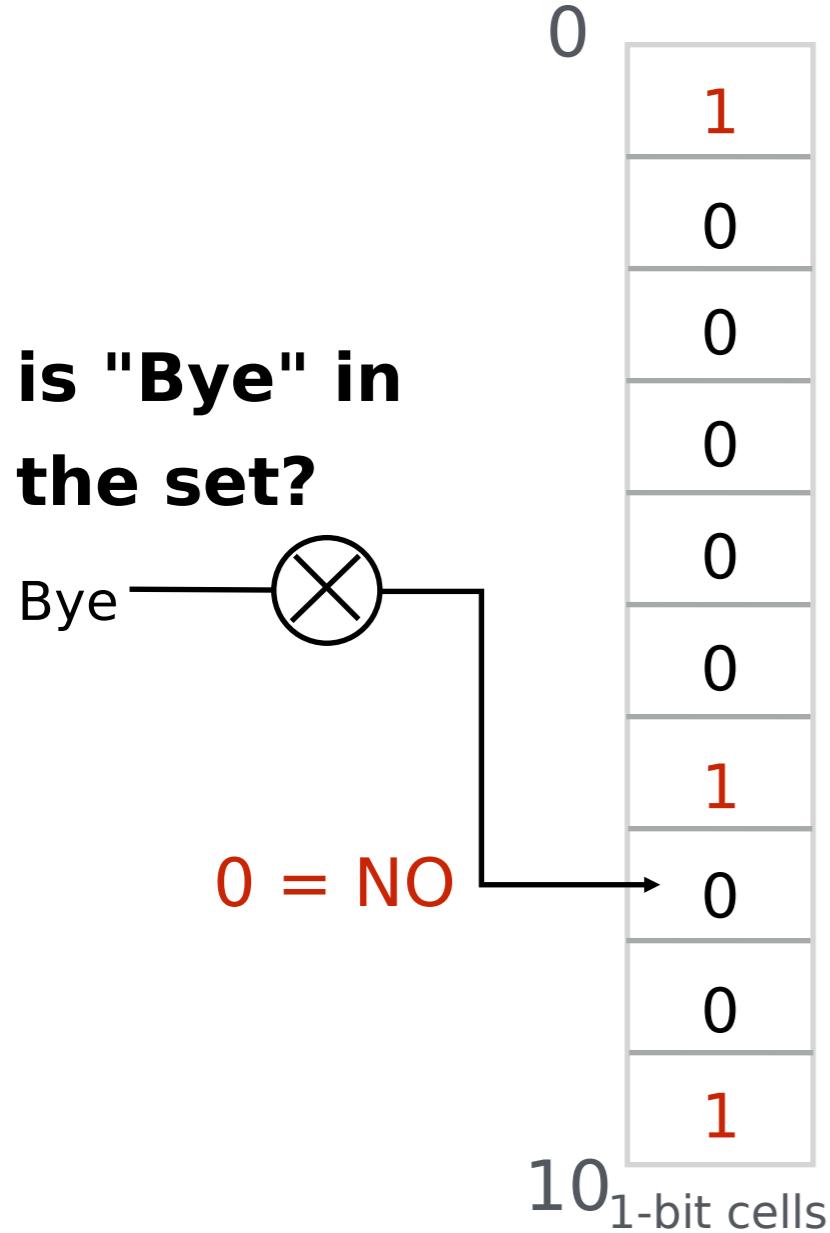
A simple approach for insertions and membership queries



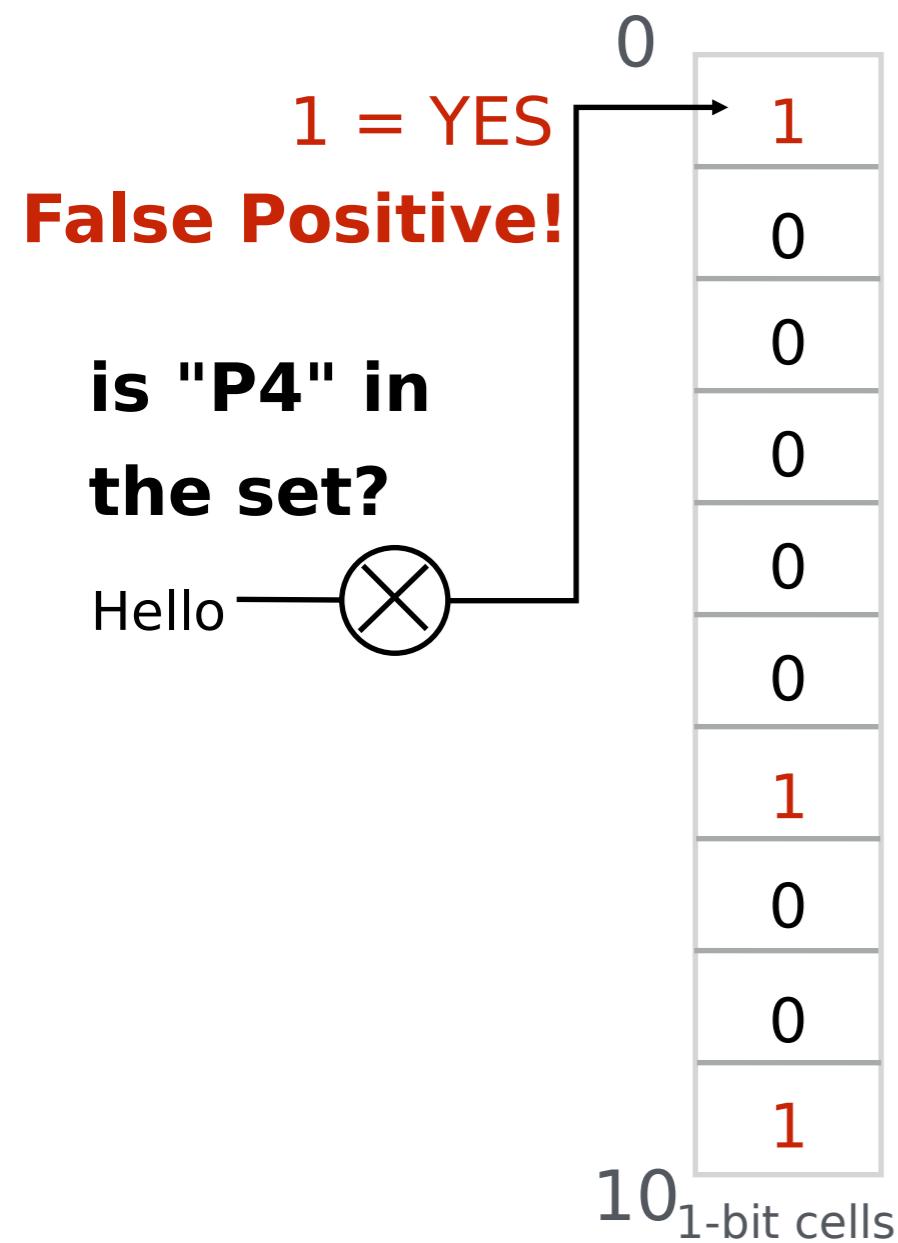
A simple approach for insertions and membership queries



A simple approach for insertions and membership queries



A simple approach for insertions and membership queries



A simple approach for insertions and membership queries

N elements and M cells

probability of an element to be mapped into a particular cell

$$\frac{1}{M}$$

probability of an element not to be mapped into a particular cell

$$1 - \frac{1}{M}$$

probability of a cell to be 0

$$(1 - \frac{1}{M})^N$$

false positive rate (FPR)

$$1 - (1 - \frac{1}{M})^N$$

false negative rate

$$0$$

A simple approach for insertions and membership queries

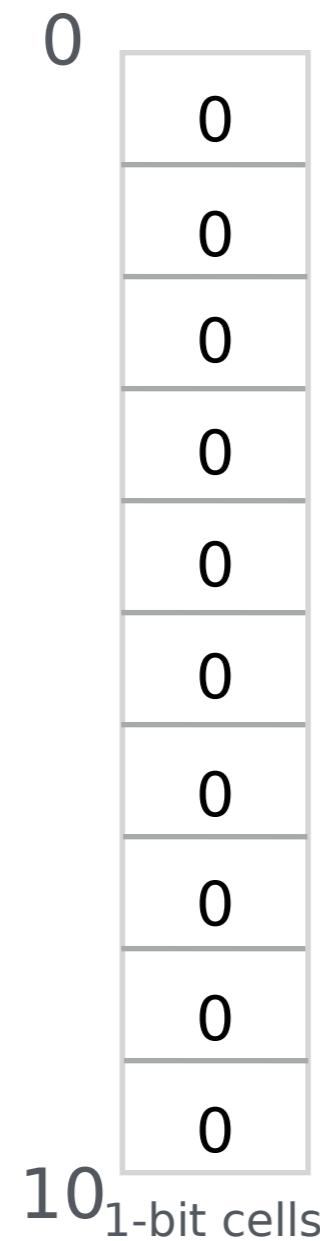
# of elements	# of cells	FPR
1000	10000	9.5%
1000	100000	1%

A simple approach for insertions and membership queries

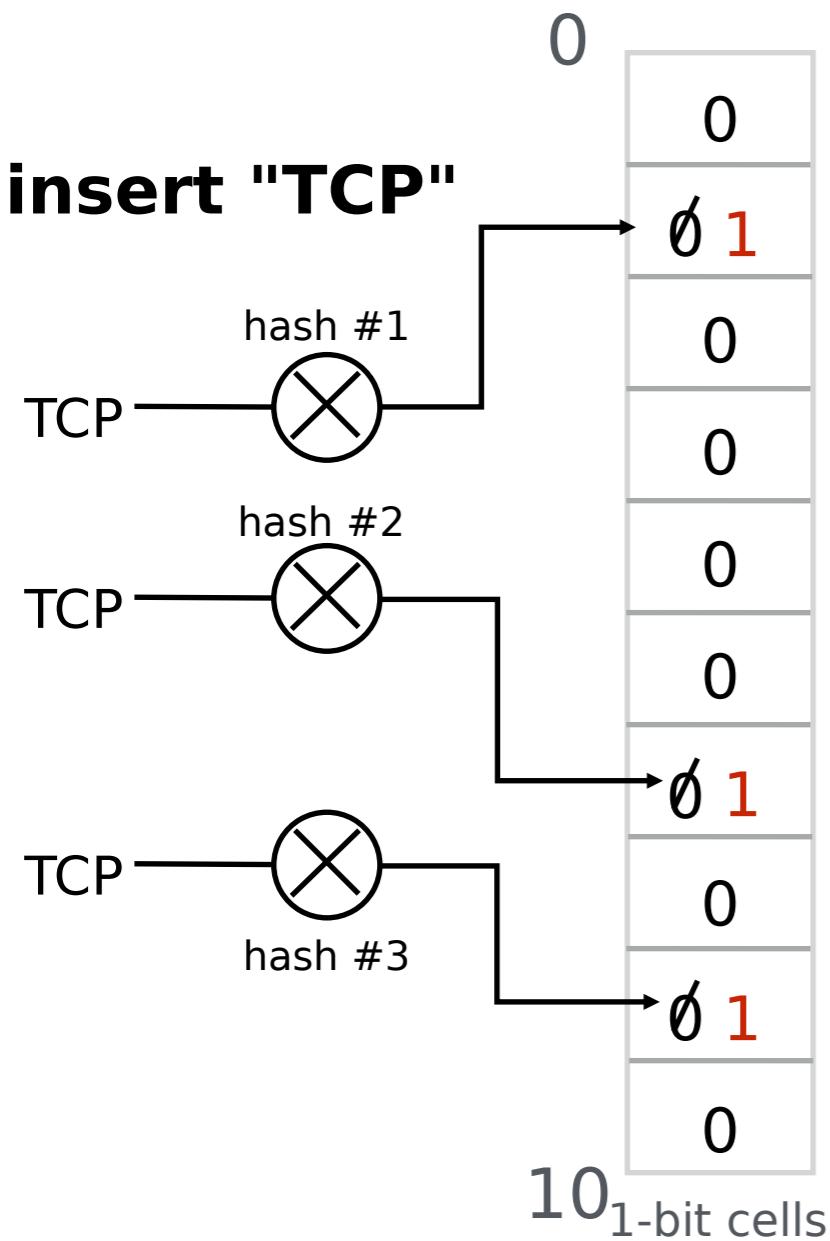
Pros: simple and only one operation per insertion or query

Con: roughly 100x more cells are required than the number of element we want to store for a 1% false positive rate

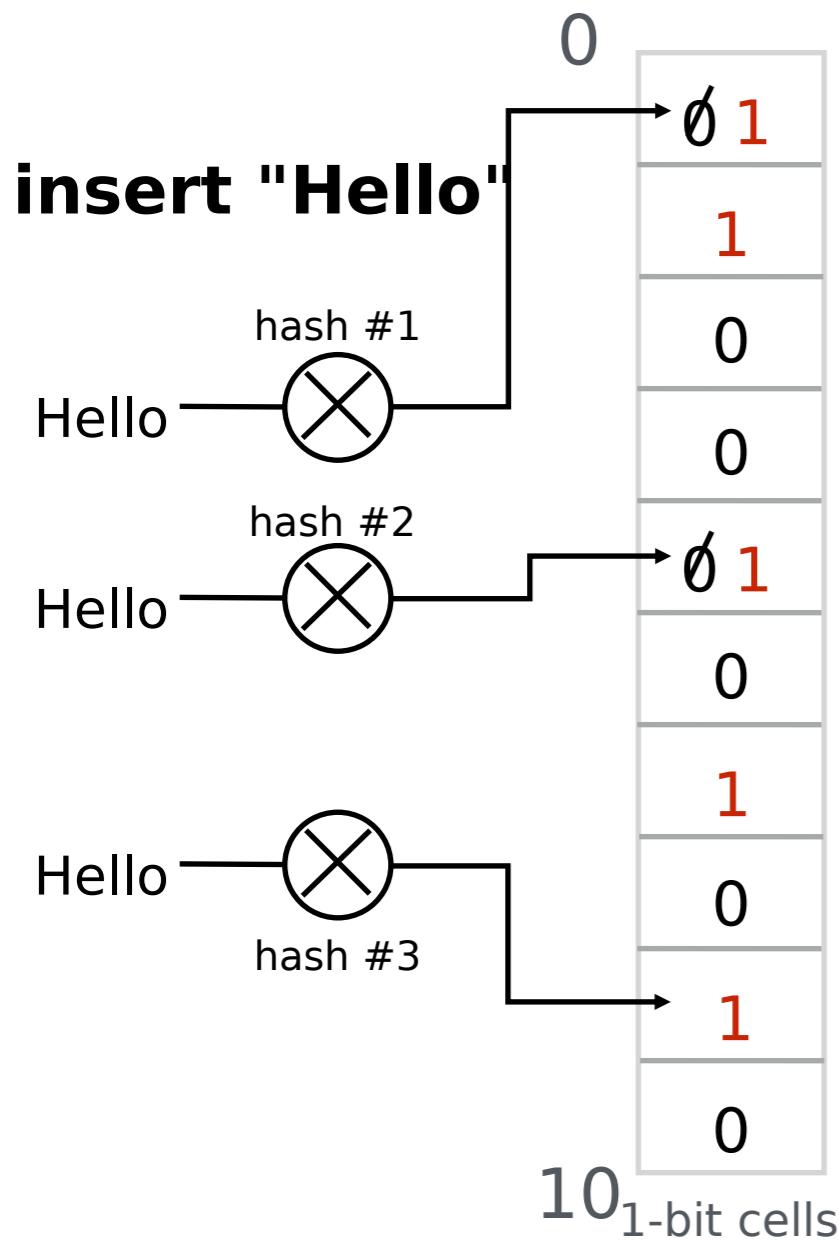
Bloom Filters: a more memory-efficient approach for insertions and membership queries



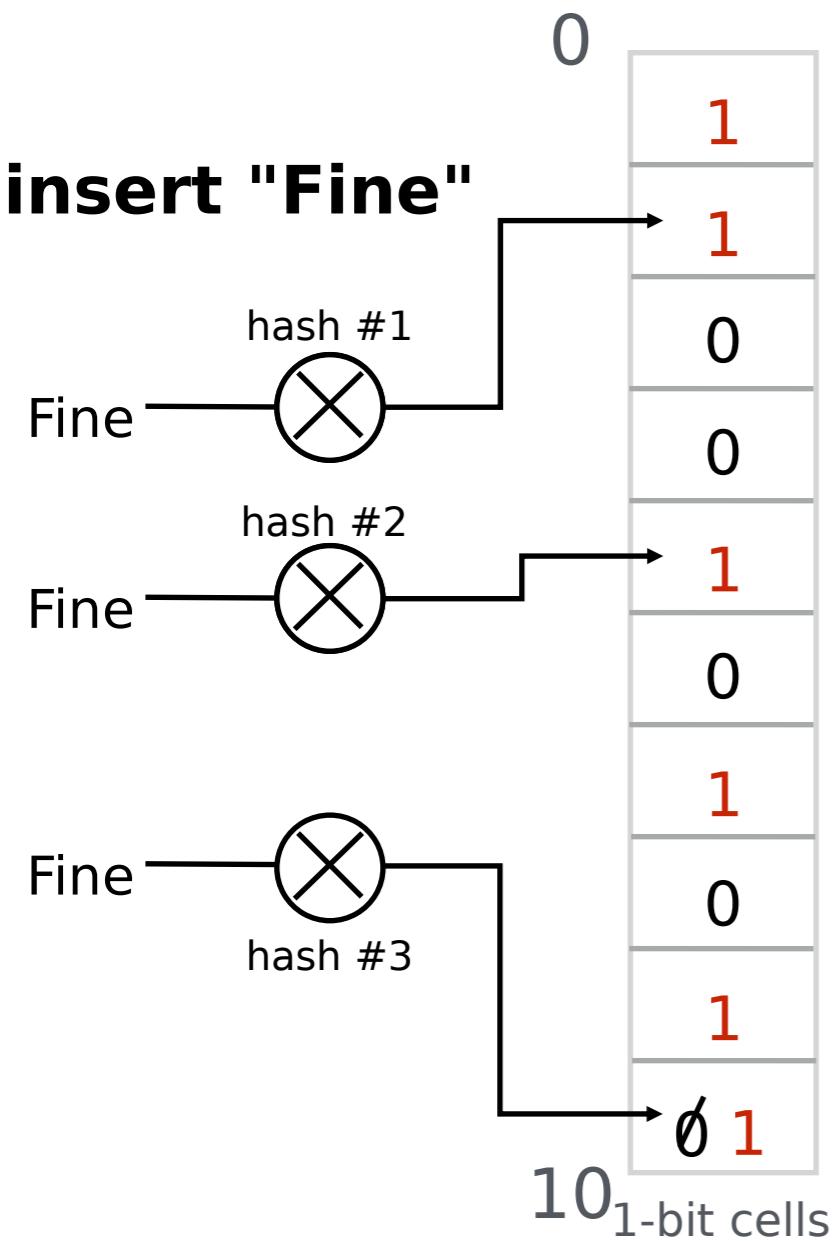
Bloom Filters: a more memory-efficient approach for **insertions** and **membership queries**



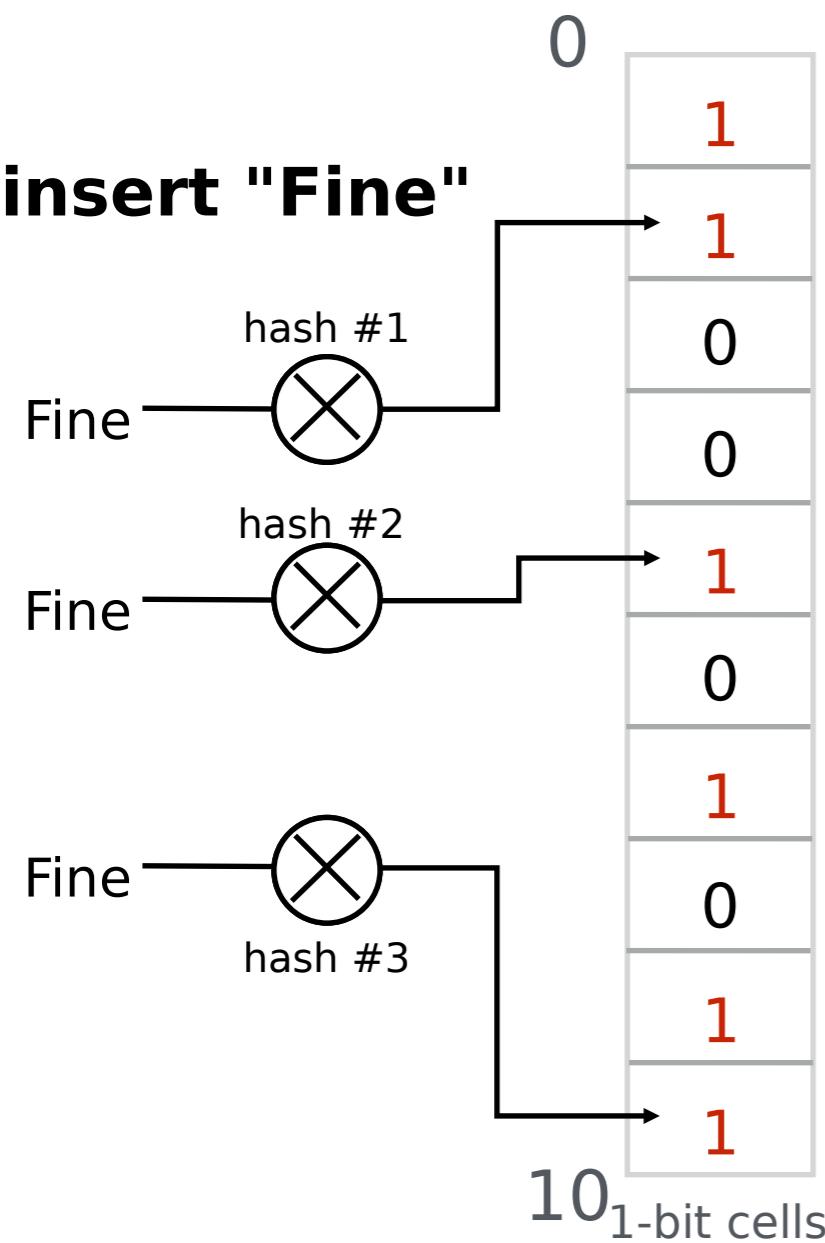
Bloom Filters: a more memory-efficient approach for **insertions** and **membership queries**



Bloom Filters: a more memory-efficient approach for **insertions** and **membership queries**



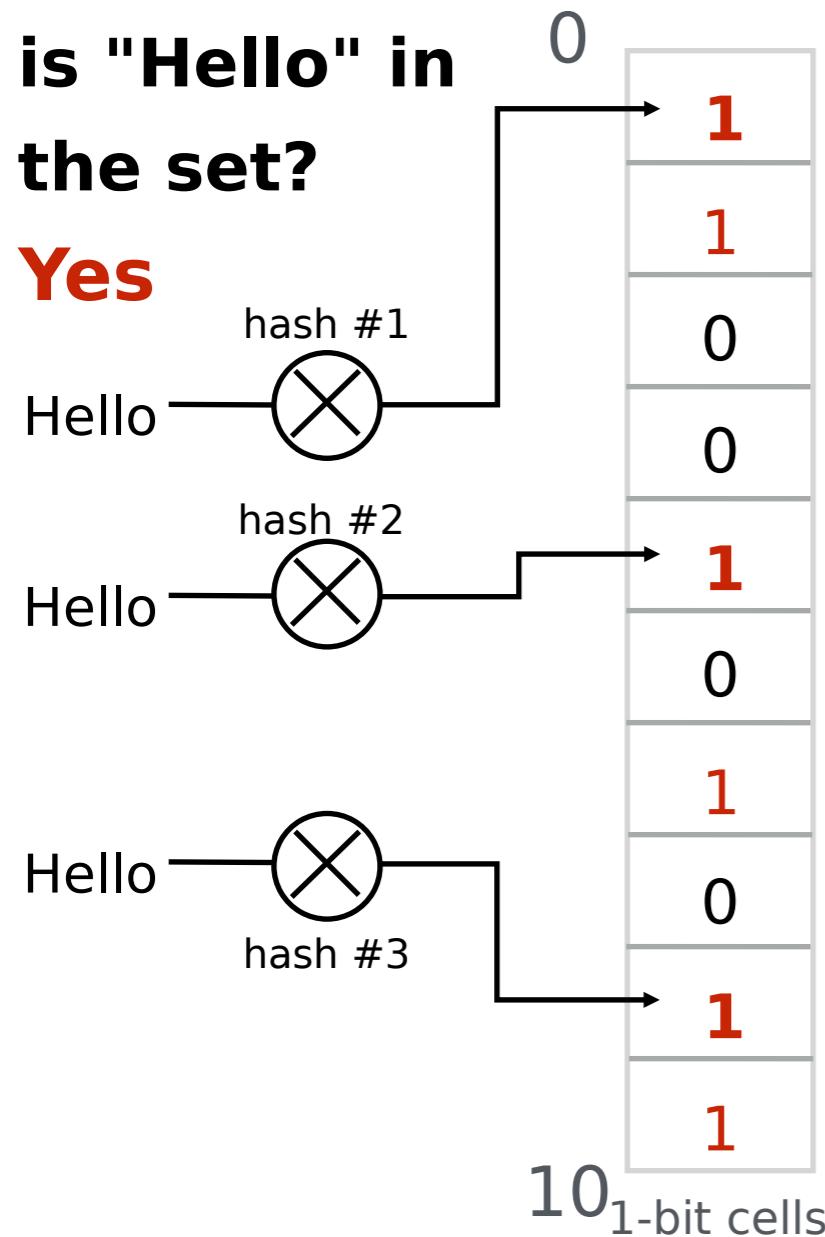
Bloom Filters: a more memory-efficient approach for **insertions** and **membership queries**



An element is considered in the set if **all** the hash values map to a cell with 1

An element is not in the set if **at least** one hash value maps to a cell with 0

Bloom Filters: a more memory-efficient approach for **insertions** and **membership queries**

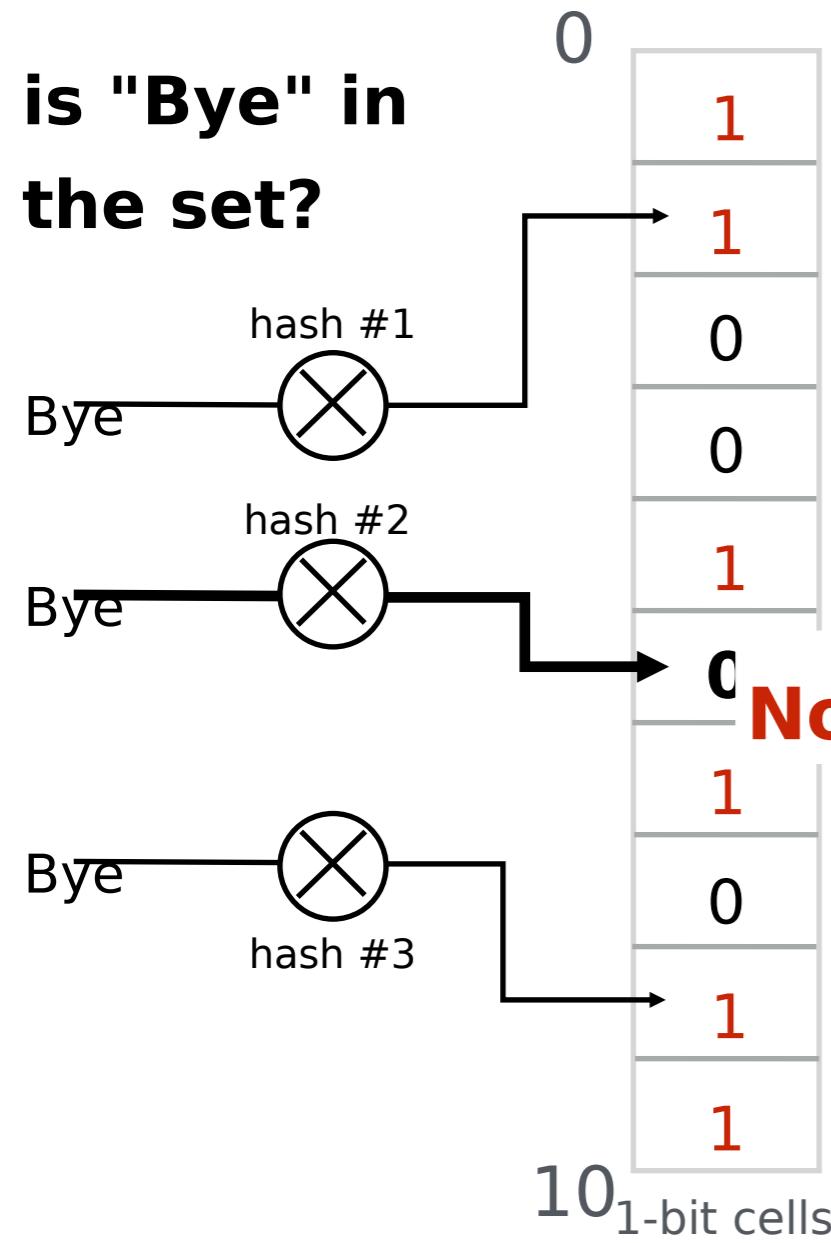


An element is considered in the set if **all** the hash values map to a cell with 1

An element is not in the set if **at least** one hash value maps to a cell with 0

Bloom Filters: a more memory-efficient approach for insertions and membership queries

**is "Bye" in
the set?**



An element is considered in the set if **all** the hash values map to a cell with 1

Do it isn't

An element is not in the set if
at least one hash value maps to
a cell with 0

Bloom Filters: a more memory-efficient approach for **insertions** and **membership queries**

False Positive!

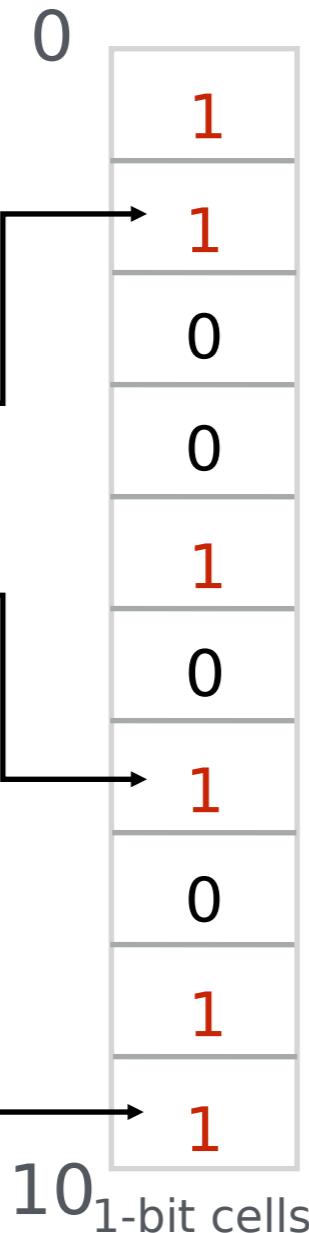
is "Fire" in
the set?

Yes

Bye → hash #1

Bye → hash #2

Bye → hash #3



An element is considered in
the set if **all** the hash values
map to a cell with 1

An element is not in the set if
at least one hash value maps to
a cell with 0

Bloom Filters: a more memory-efficient approach for **insertions** and **membership queries**

N elements, M cells and K *independent* hash functions

probability that one hash function
returns the index of a particular cell $\frac{1}{M}$

probability that one hash function does
not return the index of a particular cell $1 - \frac{1}{M}$

probability of a cell to be 0 $(1 - \frac{1}{M})^{KN}$

false positive rate $(1 - (1 - \frac{1}{M})^{KN})^K$

false negative rate 0

Bloom Filters: a more memory-efficient approach for **insertions** and **membership queries**

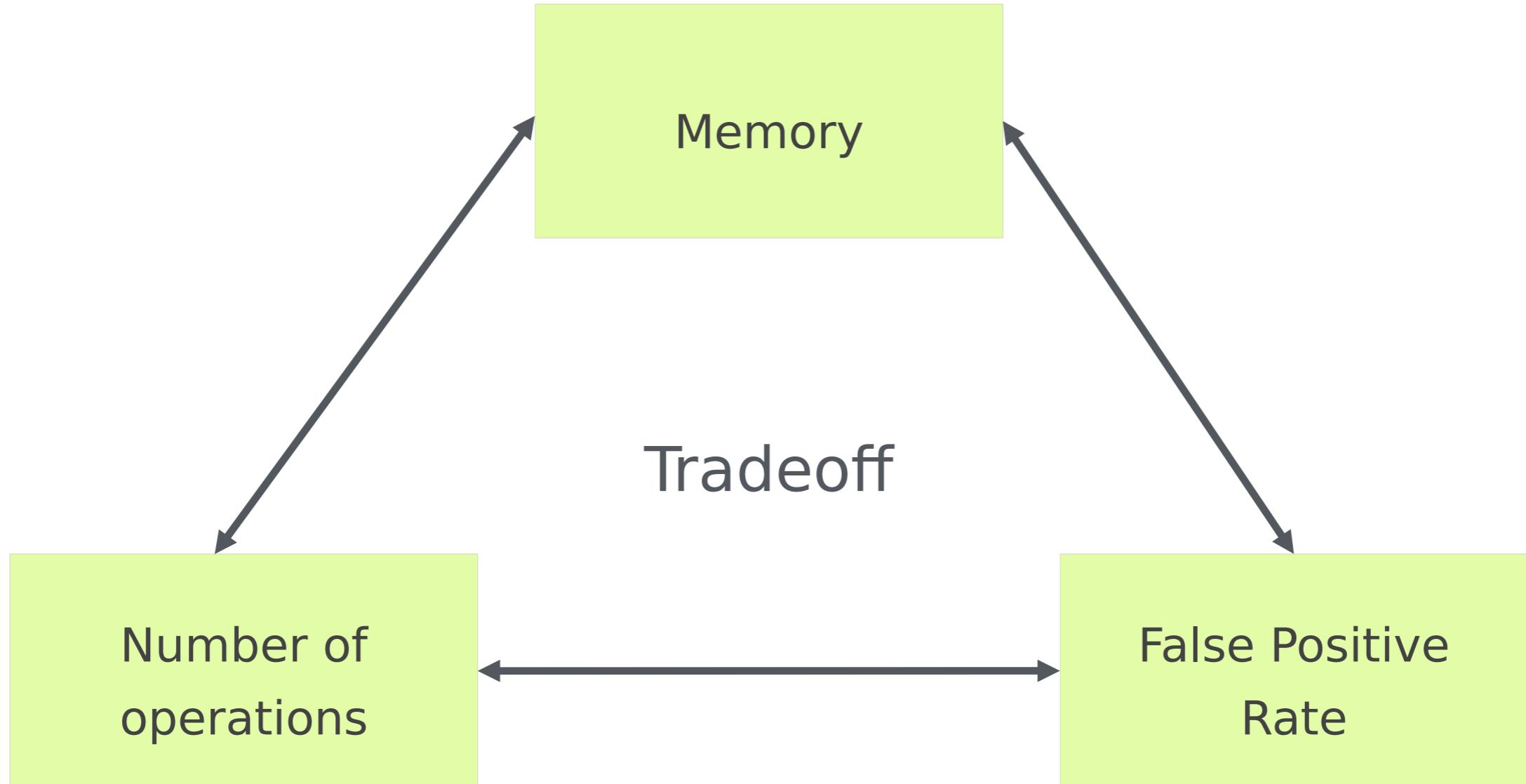
# of elements	# of cells	# hash functions	FPR
1000	10000	7	0.82%
1000	100000	7	≈ 0%

Bloom Filters: a more memory-efficient approach for **insertions** and **membership queries**

Pro: consumes roughly 10x less memory than
the simple approach

Con: Requires slightly more operations than the
simple approach (7 hashes instead of just 1)

Dimension your Bloom Filter



Dimension your Bloom Filter

N elements

M cells

K hash functions

FP false positive rate

Dimension your Bloom Filter

N elements

M cells

K hash functions

FP false positive rate

asymptotic approx.

$$FP = \left(1 - \left(1 - \frac{1}{M}\right)^{KN}\right)^K \approx \left(1 - e^{-KN/M}\right)^K$$

with calculus you can
dimension your bloom filter

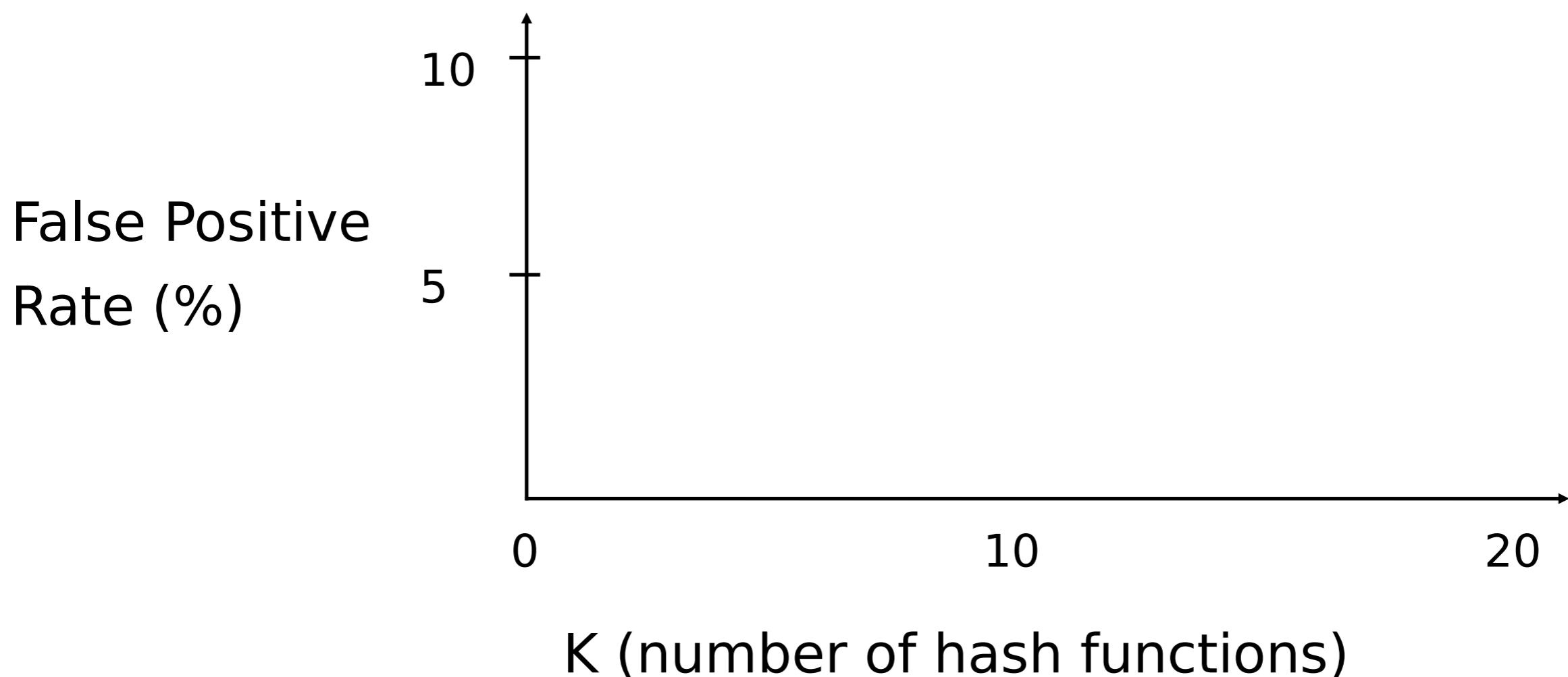
Dimension your Bloom Filter

N elements

M cells

K hash functions

FP false positive rate



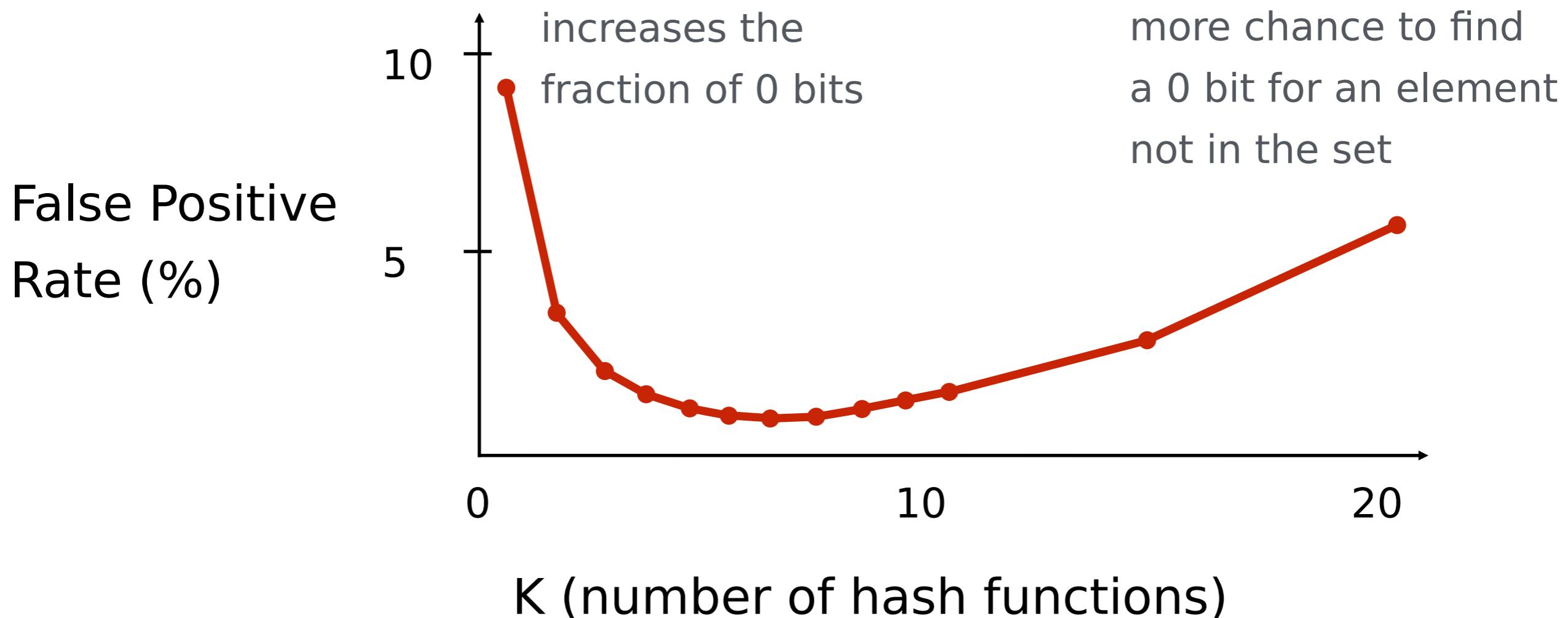
Dimension your Bloom Filter

N elements

M cells

K hash functions

FP false positive rate



Dimension your Bloom Filter

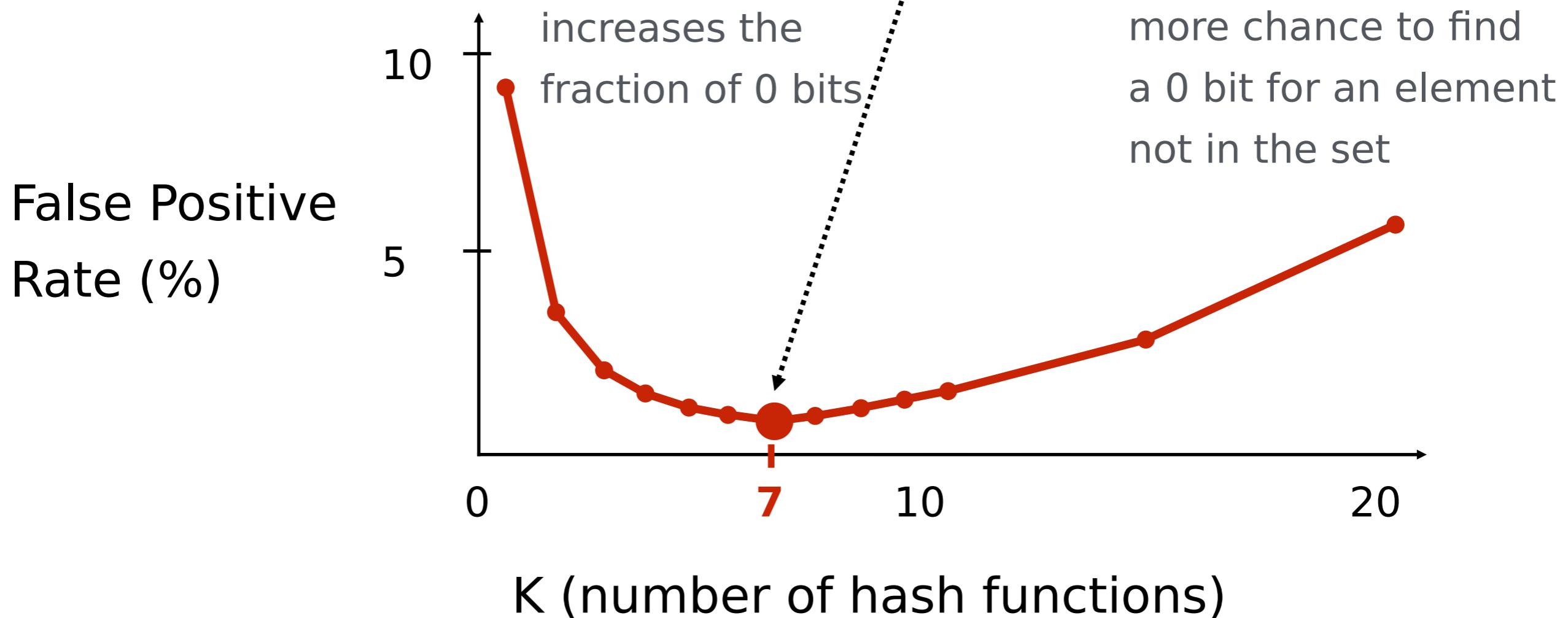
N elements

M cells

K hash functions

FP false positive rate

there is always a global minimum when $K = \ln 2 * (M/N)$ found by taking the derivative of $\approx (1 - e^{-KN/M})^K$



Implementation of a Bloom Filter in P4₁₆

You will have to use hash functions

v1model

```
enum HashAlgorithm {  
    crc32,  
    crc32_custom,  
    crc16,  
    s,  
    random,  
    identity,  
    csum16,  
    xor16  
}
```

```
extern void hash<O, T, D, M>(out O result,  
    in HashAlgorithm algo, in T base, in D data, in M  
    max);
```

more info

<https://github.com/p4lang/p4c/blob/master/p4include/v1model.p4>

Implementation of a Bloom Filter in P4₁₆

You will have to use hash functions, as well as registers

v1model

```
extern register<T> {

    register(bit<32> size);

    void read(out T result, in bit<32> index);
    void write(in bit<32> index, in T value);
}
```

Implementation of a Bloom Filter in P4₁₆ with 2 hash functions

```
control MyIngress(...) {  
    register register<bit<1>>(NB_CELLS) bloom_filter;
```

Implementation of a Bloom Filter in P4₁₆ with 2 hash functions

```
control MyIngress(...) {
    register register<bit<1>>(NB_CELLS) bloom_filter;
    apply {
        hash(meta.index1, HashAlgorithm.my_hash1, 0,
            {meta.dstPrefix, packet.ip.srcIP}, NB_CELLS);
        hash(meta.index2, HashAlgorithm.my_hash2, 0,
            {meta.dstPrefix, packet.ip.srcIP}, NB_CELLS);
```

Implementation of a Bloom Filter in P4₁₆

with 2 hash functions

```
control MyIngress(...) {
    register register<bit<1>>(NB_CELLS) bloom_filter;
    apply {
        hash(meta.index1, HashAlgorithm.my_hash1, 0,
            {meta.dstPrefix, packet.ip.srcIP}, NB_CELLS);
        hash(meta.index2, HashAlgorithm.my_hash2, 0,
            {meta.dstPrefix, packet.ip.srcIP}, NB_CELLS);

        if (meta.to_insert == 1) {
            bloom_filter.write(meta.index1, 1);
            bloom_filter.write(meta.index2, 1);
        }
        if (meta.to_query == 1) {
            bloom_filter.read(meta.query1, meta.index1);
            bloom_filter.read(meta.query2, meta.index2);

            if (meta.query1 == 0 || meta.query2 == 0) {
                meta.is_stored = 0;
            }
            else {
                meta.is_stored = 1;
            }
        }
    }
}
```

Implementation of a Bloom Filter in P4₁₆

with 2 hash functions

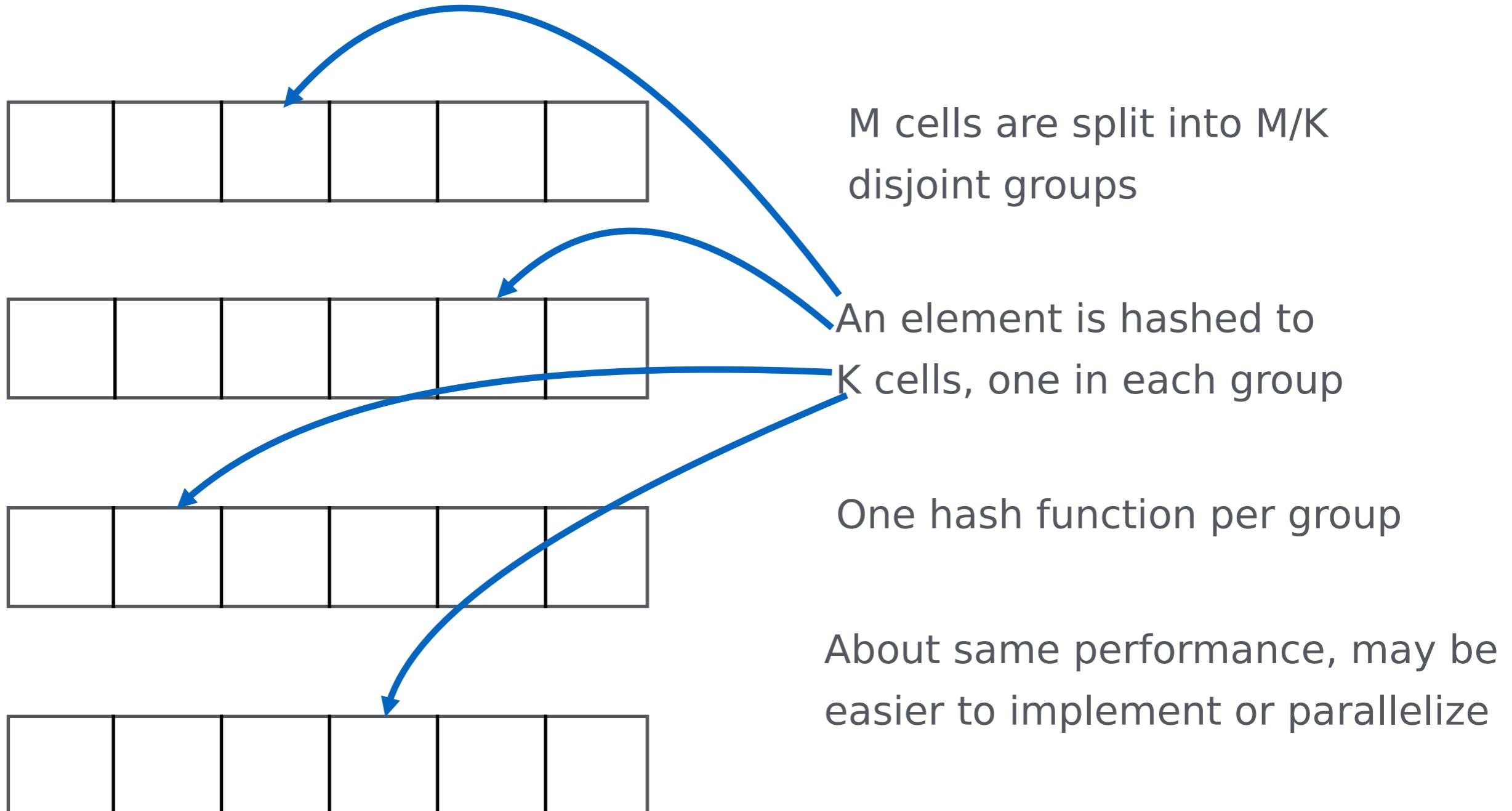
```
control MyIngress(...) {
    register register<bit<1>>(NB_CELLS) bloom_filter;
    apply {
        hash(meta.index1, HashAlgorithm.my_hash1, 0,
            {meta.dstPrefix, packet.ip.srcIP}, NB_CELLS);
        hash(meta.index2, HashAlgorithm.my_hash2, 0,
            {meta.dstPrefix, packet.ip.srcIP}, NB_CELLS);

        if (meta.to_insert == 1) {
            bloom_filter.write(meta.index1, 1);
            bloom_filter.write(meta.index2, 1);
        }
        if (meta.to_query == 1) {
            bloom_filter.read(meta.query1, meta.index1);
            bloom_filter.read(meta.query2, meta.index2);

            if (meta.query1 == 0 || meta.query2 == 0) {
                meta.is_stored = 0;
            }
            else {
                meta.is_stored = 1;
            }
        }
    }
}
```

Everything in bold
red must be adapted
for your program

Depending on the hardware limitations,
splitting the bloom filter might be required



Because deletions are not possible, the controller may need to regularly **reset** the bloom filters

Resetting a bloom filter takes some time during which it is not usable

Common trick: use two bloom filters and use one when the controller resets the other one

Bloom filters may be extended to allow deletions and to list the filter content.

If you are curious, check out the extended slides for:

- counting Bloom filters (allow deletions)
- invertible Bloom filters (allow to list content)

Bloom filters are **probabilistic data structures** for
set membership queries. For more info, see:

Space/Time Trade-offs in Hash Coding
with Allowable Errors. Burton H. Bloom. 1970.

Network Applications of Bloom Filters: A Survey.
Andrei Broder and Michael Mitzenmacher. 2004.

Invertible Bloom Lookup Tables.
Michael T. Goodrich and Michael Mitzenmacher. 2015.

FlowRadar: A Better NetFlow for Data Centers
Yuliang Li et al. NSDI 2016.



← You are looking at a stream of data (packets).
Today, I'll show you how set membership and frequency queries can be realized in P4.

PART 1

Is a certain element (e.g. ip address) in the stream?
→ Bloom filter

PART 2

How frequently does an element appear?
→ CountMin Sketch, Count Sketch, ...

part 2: counting with sketches

How frequently does an element appear?

(slides by yours truly)

We are going to look at **frequencies**,
i.e. **how often** an element occurs in a data stream.

$$\mathbf{x} = \begin{bmatrix} x_1 \\ x_2 \\ \vdots \end{bmatrix}$$

*vector of frequencies (counts)
of all **distinct elements** x_i*

We are going to look at **frequencies**,
i.e. how often an element occurs in a data stream.

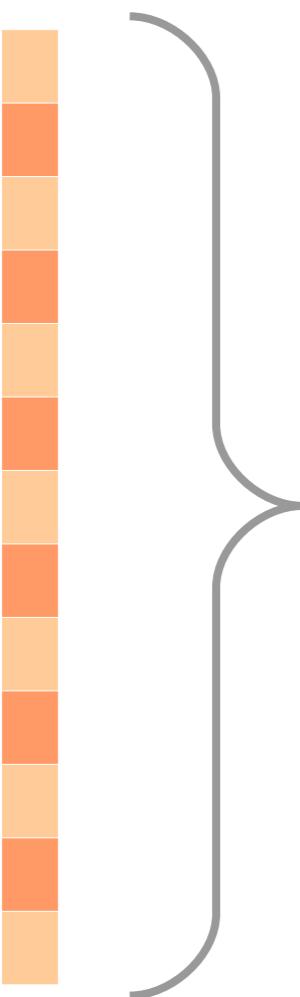
$$\mathbf{x} = \begin{bmatrix} x_1 \\ x_2 \\ \vdots \end{bmatrix}$$

*vector of frequencies (counts)
of all **distinct elements** x_i*

e.g. flows, ip addresses, ...

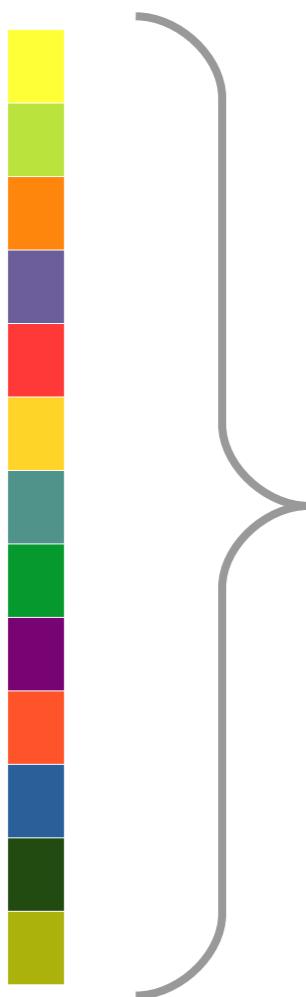
In the worst case, an algorithm providing
exact frequencies requires **linear space**.

In the worst case, an algorithm providing
exact frequencies requires linear space.



Data Stream
n elements in total

In the worst case, an algorithm providing exact frequencies requires linear space.



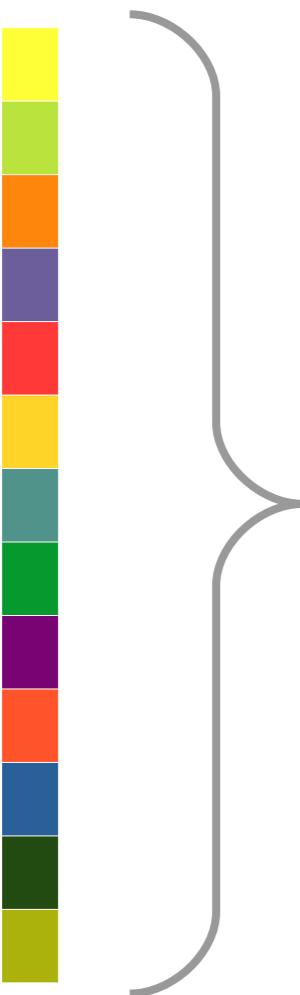
Data Stream

n elements in total

→ **n distinct elements**

(in the worst case)

In the worst case, an algorithm providing
exact frequencies requires linear space.



Data Stream

n elements in total

→ **n distinct elements**

(in the worst case)

→ **n counters required? :(**

Probabilistic datastructures can help again!

Bloom Filters

quickly “filter” only those elements that might be in the set

More efficient by allowing false positives.

Probabilistic datastructures can help again!

Bloom Filters

quickly “filter” only those elements that might be in the set

More efficient by allowing false positives.

Sketches

provide approximate frequencies of elements in a data stream.

More efficient by allowing mis-counting.

A CountMin sketch is designed to have provable L1 error bounds for frequency queries.

A CountMin sketch is designed to have provable L1 error bounds for frequency queries.



Notation reminder:
vector of frequencies (counts)
of all distinct elements x_i

$$x = \begin{bmatrix} x_1 \\ x_2 \\ \vdots \end{bmatrix}$$

A CountMin sketch is designed to have provable
L1 error bounds for frequency queries.

$$Pr [\hat{x}_i - x_i \geq \varepsilon \|x\|_1] \leq \delta$$

estimated frequency *true frequency* *sum of frequencies*

The estimation error exceeds $\varepsilon \|x\|_1$
with a probability smaller than δ

$$Pr [\hat{x}_i - x_i \geq \varepsilon \|\mathbf{x}\|_1] \leq \delta$$

estimated frequency *true frequency* *sum of frequencies*

relative to L1 norm

The estimation error exceeds $\varepsilon \|\mathbf{x}\|_1$ with a probability smaller than δ

$$Pr [\hat{x}_i - x_i \geq \varepsilon \|x\|_1] \leq \delta$$

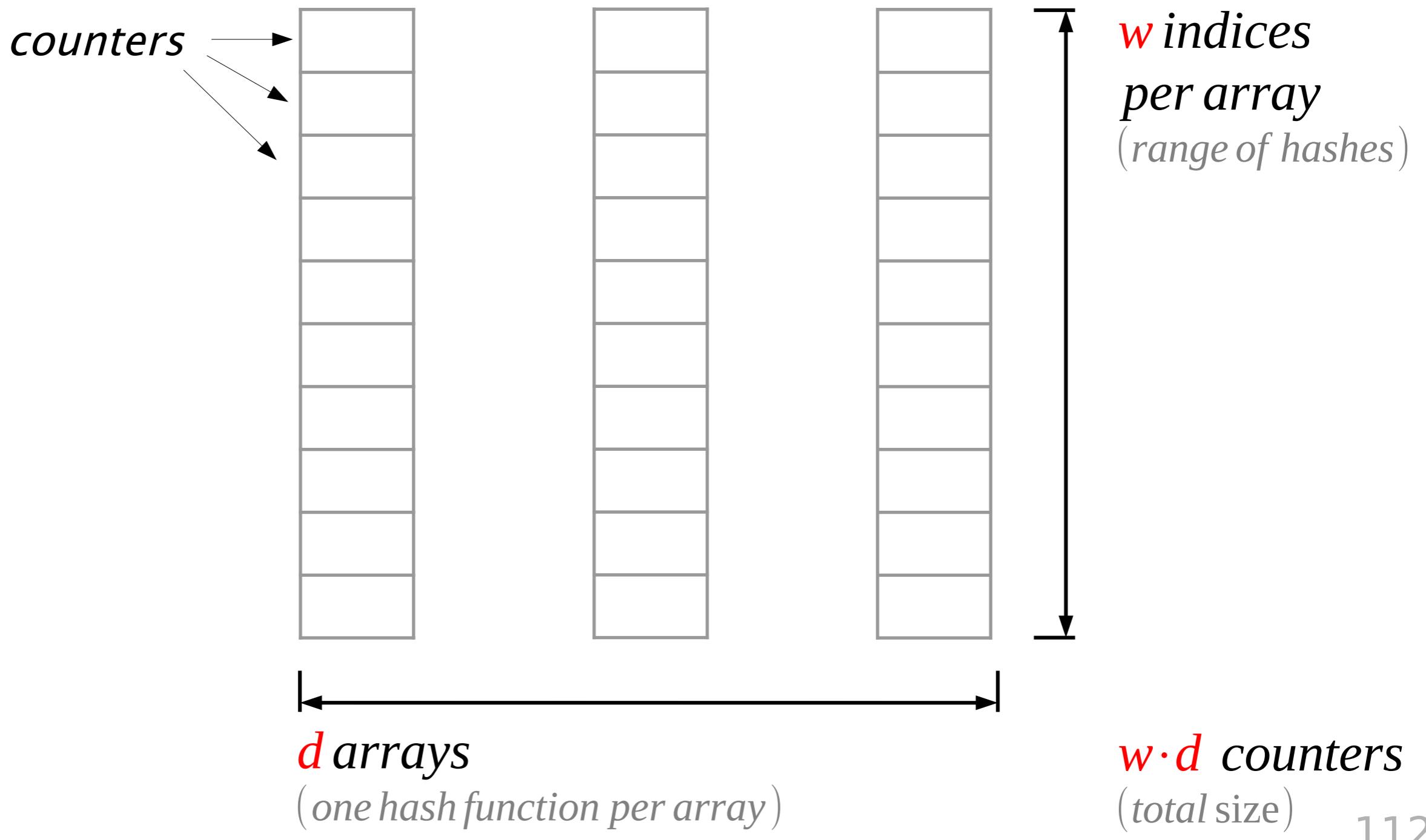
estimated frequency *true frequency* *sum of frequencies*

Let $\|x\|_1 = 10000$, $\varepsilon = 0.01$, $\delta = 0.05$

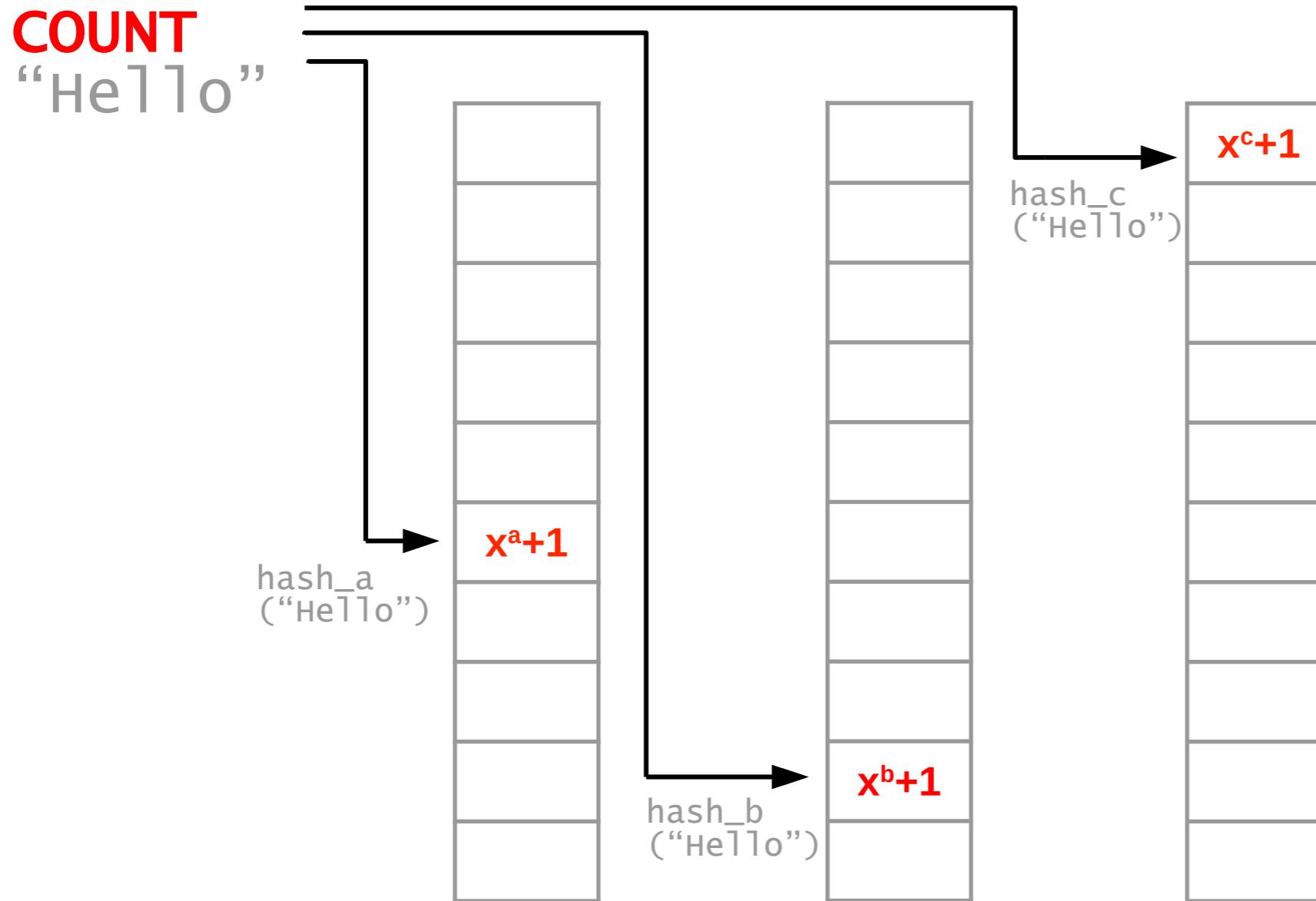
*Then, after counting 10000 elements in total,
 the probability for any estimate to be
 off by more than 100 is less than 5%.*

A CountMin sketch is **designed** to have provable L1 error bounds for frequency queries.

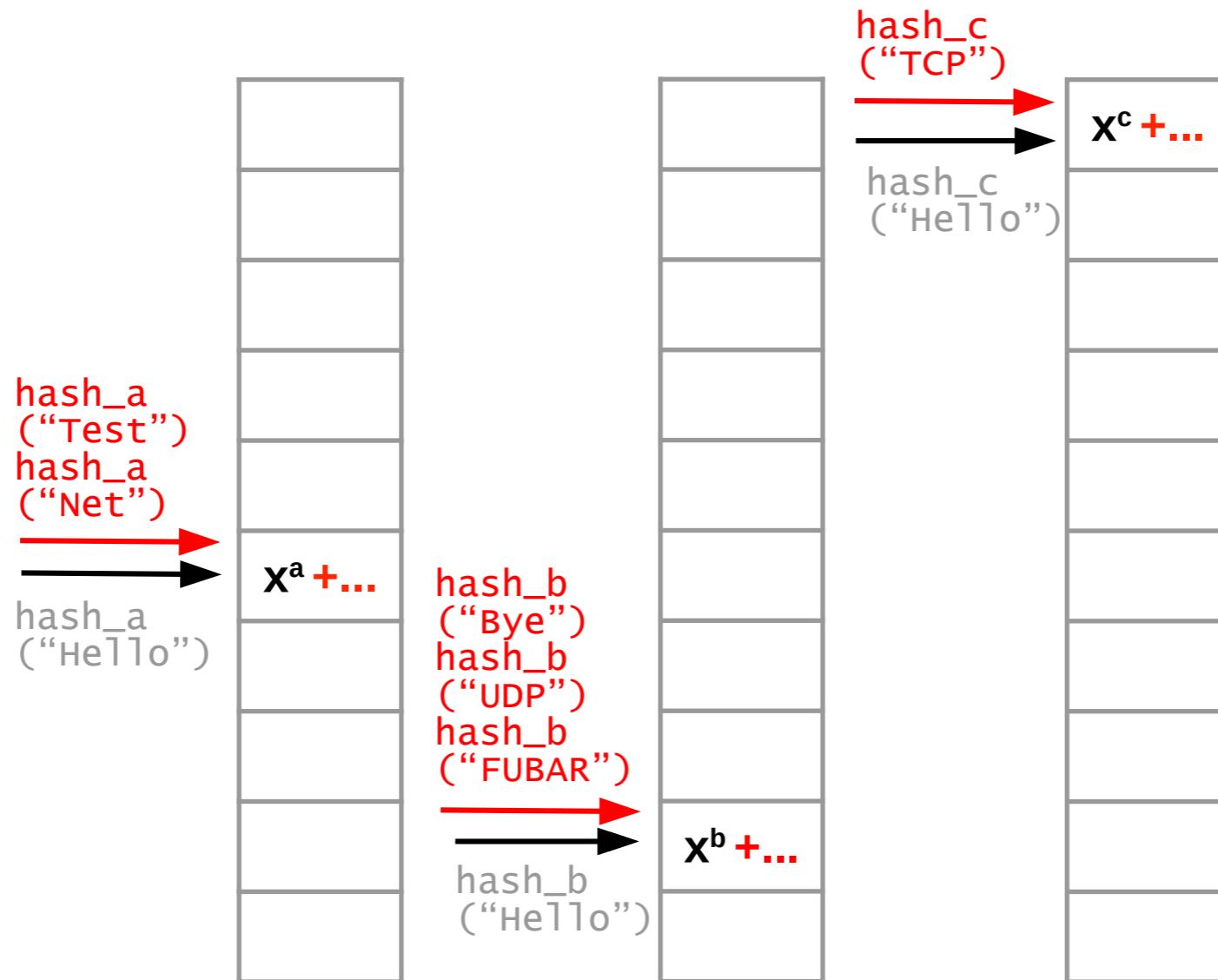
A CountMin Sketch uses multiple arrays and hashes.



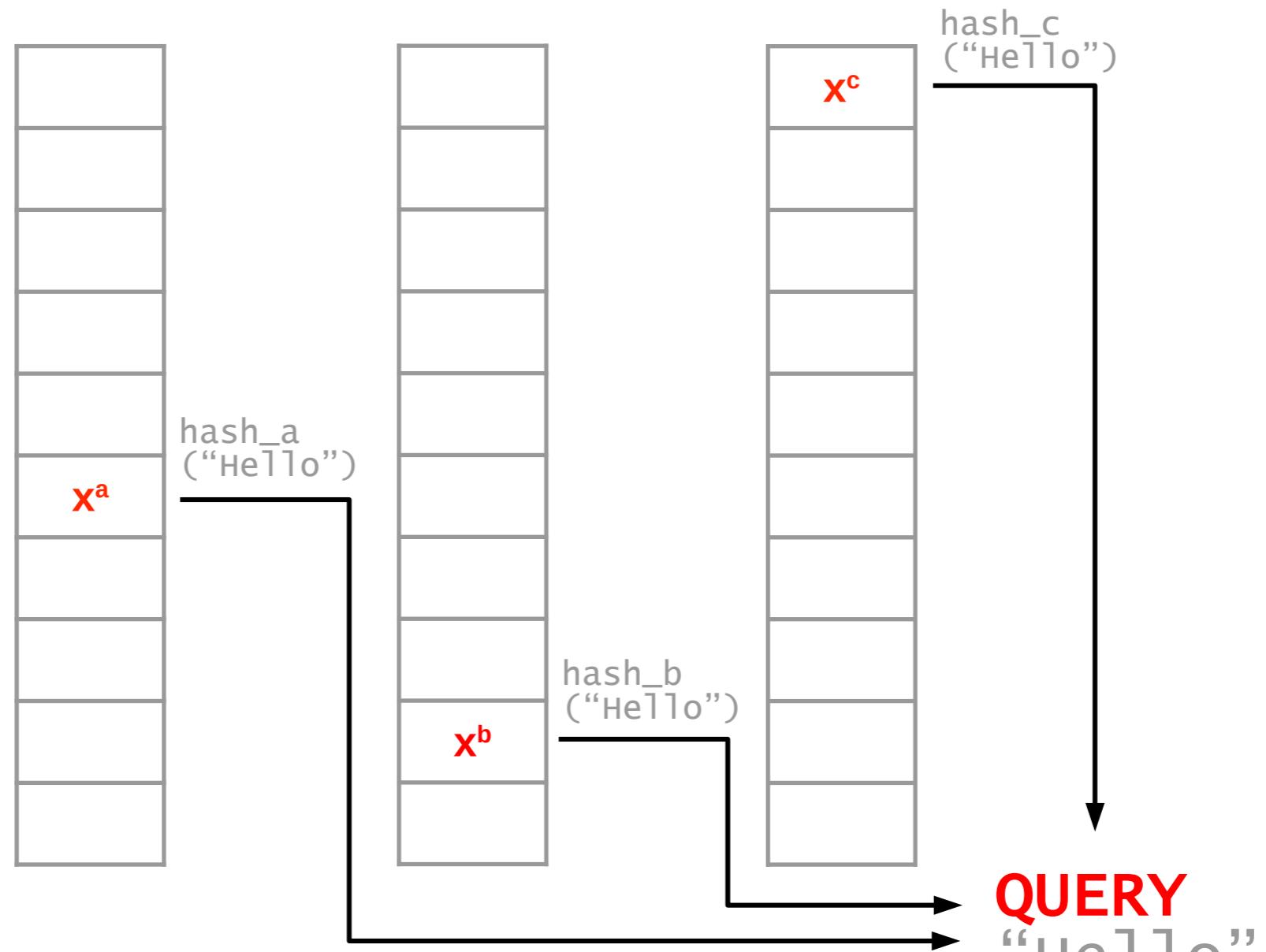
To count, increment all hash-indexed fields by 1.



Hash collisions cause over-counting.



Returning the **minimum value** minimizes the error.



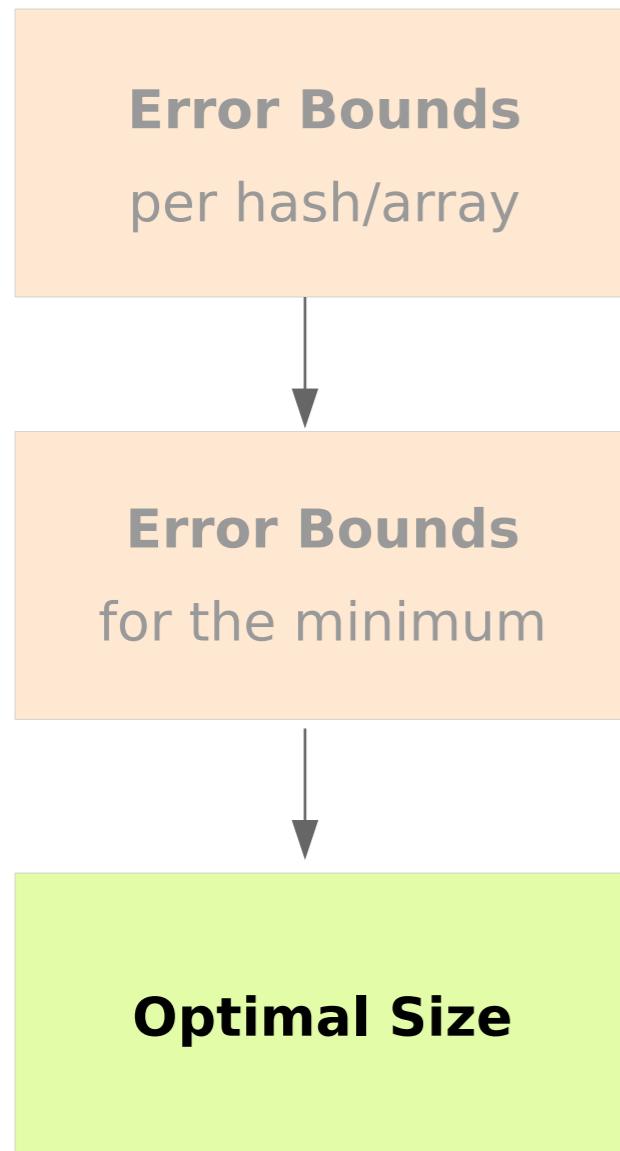
A CountMin sketch is designed to have **provable L1 error bounds** for frequency queries.

$$Pr [\hat{x}_i - x_i \geq \varepsilon \|x\|_1] \leq \delta$$

estimated frequency *true frequency* *sum of frequencies*

A CountMin sketch recipe

(see extended slides for derivation)



Given ε, δ , choosing

$$w = \left\lceil \frac{e}{\varepsilon} \right\rceil \quad (\text{hash range})$$

$$d = \left\lceil \ln \frac{1}{\delta} \right\rceil \quad (\#\text{hashes})$$

requires the **minimum number of counters** s.t. the CountMin Sketch can guarantee that

$$\hat{x}_i - x_i \geq \varepsilon \|x\|_1$$

with a probability less than δ

A CountMin sketch is designed to have provable L1 error bounds for frequency queries.

CountMin sketch recipe

Choose $d = \left\lceil \ln \frac{1}{\delta} \right\rceil, w = \left\lceil \frac{e}{\varepsilon} \right\rceil$

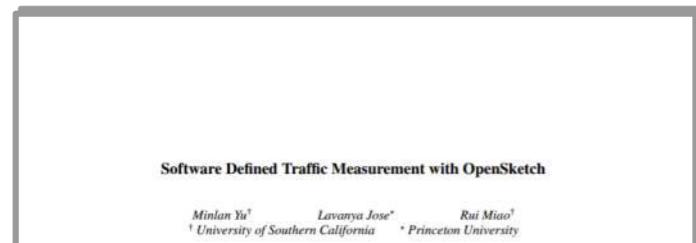
Then $\hat{x}_i - x_i \geq \varepsilon \|x\|_1$ with a probability less than δ

Sketches are the new black

OpenSketch

NSDI '13

[source]



Abstract

Most network management tasks in software-defined networks (SDN) involve two stages: measurement and control. While many efforts have been focused on network control APIs for SDN, little attention goes into measurement. The key challenge of designing a new measurement API is to strike a careful balance between generality (supporting a wide variety of measurement tasks) and efficiency (enabling high link speed and low cost). We propose a software defined traffic measurement system called OpenSketch which separates the measurement data plane from the control plane. In the data plane, OpenSketch provides a simple three-stage pipeline (hashing, filtering, and counting), which can be implemented with commodity switch components and support many measurement tasks. In the control plane, OpenSketch provides a measurement library that automatically configures the pipeline and allocates resources for different measurement tasks. Our evaluations of real-world packet traces, our prototype on NetFPGA, and the implementation of five measurement tasks on top of OpenSketch, demonstrate that OpenSketch is general, efficient and easily programmable.

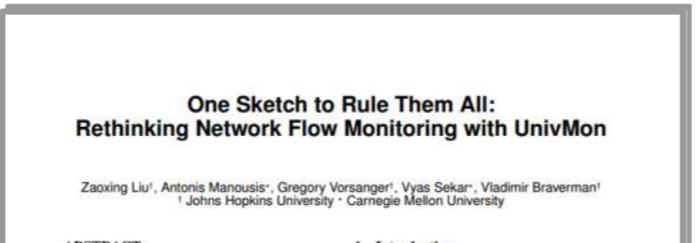
1 Introduction

Recent advances in software-defined networking (SDN) have significantly improved network management. Network management involves two important stages: (1) measuring the network in real time (e.g., identifying traffic anomalies or large traffic aggregates) and then (2) adjusting the control of the network accordingly (e.g., routing, access control, and rate limiting). While there have been many efforts on designing the right APIs for network control (e.g., OpenFlow [29], ForCES [1], rule-based forwarding [33], etc.), little thought has gone into designing the right APIs for measurement. Since con-

UnivMon

SIGCOMM '16

[source]



ABSTRACT

Network management requires accurate estimates of metrics for many applications including traffic engineering [11, 32], attack and anomaly detection [49], and forensic analysis [46]. Each such management task requires accurate and timely statistics of different application-level metrics of interest, e.g., the flow size distribution [37], heavy hitters [10], entropy measures [38, 50], or detecting changes in traffic patterns [44]. At a high level, there are two classes of techniques to estimate these metrics of interest. The first class of approaches relies on generic flow monitoring, typically with some form of packet sampling (e.g., NetFlow [25]). While generic flow monitoring is good for coarse-grained visibility, prior work has shown that it provides low accuracy for more fine-grained metrics [30, 31, 43]. These well-known limitations of sampling-based monitoring have driven the development of alternative sketching or streaming algorithms. Here, custom online algorithms and data structures are designed for specific metrics of interest that can yield provable resource-accuracy trade-offs (e.g., [17, 18, 20, 31, 36, 38, 43]).

As an alternative, many sketch-based streaming algorithms have been proposed in the theoretical research community [7, 12, 46, 8, 20, 47], which provide efficient measurement support for individual management tasks. However, these algorithms are not deployed in practice because of their lack of generality: Each of these algorithms answers just one question or produces just one statistic (e.g., the unique number of destinations), so it is too expensive for vendors to build new hardware to support each function. For example, the Space-Saving heavy hitter detection algorithm [8] maintains a hash table of items and counts, and requires customized operations such as keeping a pointer to the item with minimum counts and replacing the minimum-count entry with a

CCS Concepts

• Networks → Network monitoring; Network measurement;

Keywords

Flow Monitoring, Sketching, Streaming Algorithms
Permission to make digital or hard copies of all or part of this work for personal use or internal group use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than ACM must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from www.acm.org.

© 2016 ACM. ISBN 978-1-4503-4567-4/16/08...\$15.00
<https://doi.org/10.1145/2953482.2954906>

SketchLearn

SIGCOMM '18

[source]

SketchLearn: Relieving User Burdens in Approximate Measurement with Automated Statistical Inference

Qun Huang[†], Patrick P. C. Lee[‡], and Yungang Bao[†]

[†]State Key Lab of Computer Architecture, Institute of Computing Technology, Chinese Academy of Sciences

[‡]Department of Computer Science and Engineering, The Chinese University of Hong Kong

ABSTRACT

Network measurement is challenged to fulfill stringent resource requirements in the face of massive network traffic. While approximate measurement can trade accuracy for resource savings, it demands manual efforts to configure the right resource-accuracy trade-offs in real deployment. Such user burdens are caused by how existing approximate measurement approaches inherently deal with resource conflicts when tracking massive network traffic with limited resources. In particular, they tightly couple resource configurations with accuracy parameters, so as to provision sufficient resources to bound the measurement errors. We design SketchLearn, a novel sketch-based measurement framework that leverages machine learning and statistical properties to eliminate conflicting traffic components. We prototype SketchLearn on OpenSwitch and P4, and our testbed experiments and stress-test simulation show that SketchLearn accurately and automatically monitors various traffic statistics and effectively supports network-wide measurement with limited resources.

CCS CONCEPTS

• Networks → Network measurement;

KEYWORDS

Sketch; Network measurement

ACM Reference Format:
Qun Huang, Patrick P. C. Lee, and Yungang Bao. 2018. SketchLearn: Relieving User Burdens in Approximate Measurement with Automated Statistical Inference. In *SIGCOMM '18: ACM SIGCOMM* (Budapest, Hungary, August 20–25, 2018), 17 pages. <https://doi.org/10.1145/3230543.3230559>

Permission to make digital or hard copies of all or part of this work for personal use or internal group use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than ACM must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from www.acm.org.

© 2018 Association for Computing Machinery.
ACM ISBN 978-1-4503-5567-4/18/08...\$15.00
<https://doi.org/10.1145/3230543.3230559>

1 INTRODUCTION

Network measurement is indispensable to modern network management in clouds and data centers. Administrators measure a variety of traffic statistics, such as per-flow frequency, to infer the key behaviors or any unexpected patterns in operational networks. They use the measured traffic statistics to form the basis of management operations such as traffic engineering, performance diagnosis, and intrusion prevention. Unfortunately, measuring traffic statistics is non-trivial in the face of massive network traffic and large-scale network deployment. Error-free measurement requires per-flow tracking [15], yet today's data center networks can have thousands of concurrent flows in a very small period from 50ms [2] down to even 5ms [56]. This would require tremendous overhead for per-flow per-flow tracking.

In view of the resource constraints, researchers in the literature leverage approximation techniques to trade between resource usage and measurement accuracy. Examples include sampling [9, 37, 64], top-k counting [5, 43, 44, 46], and sketch-based approaches [18, 33, 40, 42, 58], which we collectively refer to as *approximate measurement* approaches. Their idea is to construct compact sub-linear data structures to record traffic statistics, backed by theoretical guarantees on how to achieve accurate measurement with limited resources. Approximate measurement has formed building blocks for many state-of-the-art network-wide measurement systems (e.g., [32, 48, 55, 60, 62, 67]), and is also adopted in production data centers [31, 68].

Although theoretically sound, existing approximate measurement approaches are inconvenient for use. In such approaches, massive servers or traffic competes for the limited resources, thereby introducing measurement errors due to *resource conflicts* (e.g., multiple flows are mapped to the same counter in sketch-based measurement). To mitigate errors, sufficient resources must be provisioned in approximate measurement based on its theoretical guarantees. Thus, *there exists a tight binding between resource configurations and accuracy parameters*. Such tight binding leads to several practical limitations (see §2.2 for details): (i) administrators need

Sketches are the new black

LightGuardian
NSDI '21

[source]

LightGuardian: A Full-Visibility, Lightweight, In-band Telemetry System Using Sketchlets

Yikai Zhao[†], Kaicheng Yang[†], Zirui Liu[†], Tong Yang^{†,§}, Li Chen[‡], Shiyi Liu[†], Naiqian Zheng[†], Ruixin Wang[†], Hanbo Wu[†], Yi Wang^{†,§}, Nicholas Zhang[†]

[†]Department of Computer Science, Peking University, China
[‡]Peng Cheng Laboratory, Shenzhen, China [§]Huawei Theory Lab, China

^{†,§}Southern University of Science and Technology

Abstract

Network traffic measurement is central to successful network operations, especially for today's hyper-scale networks. Although existing works have made great contributions, they fail to achieve the following three criteria simultaneously: 1) **full-visibility**, which refers to the ability to acquire any desired per-flow level information for all flows; 2) **low overhead** in terms of computation, memory, and bandwidth; and 3) **robustness**. We design LightGuardian to meet these three criteria. Our key innovation is a novel constant-sized data structure, called *sketchlet*, which can be embedded in packet headers. Specifically, we design a novel SubMax sketch to accurately capture flow-level information. SubMax can be divided into sketchlets, which are carried in-band by passing packets to the end-hosts for aggregation, reconstruction, and analysis. We have fully implemented a LightGuardian prototype on a testbed with 10 programmable switches and 8 end-hosts in a FatTree topology, and conduct extensive experiments and evaluations. Experimental results show that LightGuardian can obtain per-flow per-hop flow-level information within 1.0 ~ 1.5 seconds with consistently low overhead, using only 0.07% total bandwidth capacity of the network. We believe LightGuardian is the first system to collect per-flow per-hop information for all flows in the network with negligible overhead.

1 Introduction

Network traffic measurement is central to successful network operations, especially for today's hyper-scale networks with more than 10⁹ devices [1–6]. Meanwhile, at end-hosts, knowing the traffic information in the core of the network can also benefit application performance [7–9]. To infer application performance and user experience, the community consensus is to measure at flow-level granularity. Thus, an ideal measurement system is expected to achieve: (1) **full-visibility**, which we define as the ability to acquire any desired per-

Nearly-Zero-Error
NSDI '21

[source]

Toward Nearly-Zero-Error Sketching via Compressive Sensing

Qun Huang^{1,2} Siyuan Sheng³ Xiang Chen^{1,2,4} Yungang Bao⁵ Rui Zhang⁵ Yanwei Xu⁵ Gong Zhang⁵
¹Peking University ²Pengcheng Lab ³Institute of Computing Technology, CAS
⁴Fuzhou University ⁵Huawei Theory Department

Abstract

Sketch algorithms have been extensively studied in the area of network measurement, given their limited resource usage and theoretically bounded errors. However, error bounds provided by existing algorithms remain too coarse-grained: in practice, only a small number of flows (e.g., heavy hitters) actually benefit from the bounds, while the remaining flows still suffer from serious errors. In this paper, we aim to design a nearly-zero-error sketch that achieves negligible per-flow error for almost all flows. We base our study on a technique named compressive sensing. We exploit compressive sensing in two aspects. First, we incorporate the near-perfect recovery of compressive sensing to boost sketch accuracy. Second, we propose a new sketching mechanism based on a novel methodology to analyze various choices of sketch algorithms. Guided by the analysis, we propose two sketch algorithms that seamlessly embrace compressive sensing to reach nearly zero errors. We implement our algorithms in OpenVSwitch and P4. Experimental results show that the two algorithms incur less than 0.1% per-flow error for more than 99.72% flows, while preserving the resource efficiency of sketch algorithms. The efficiency demonstrates the power of our new methodology for sketch analysis and design.

1 Introduction

Sketch algorithms have been widely adopted in flow-level monitoring. They maintain compact data structures that sacrifice a small portion of accuracy to be readily deployable in commodity network devices. Given their limited overheads and provable high accuracy, numerous sketch algorithms are designed to monitor various flow statistics, such as per-flow counts [49], heavy hitters [19, 25], denoted-service vice counts [26, 40], and so on [40]. These flow statistics form essential building blocks for network management.

Despite the sound theoretical bounds on the errors, existing algorithms remain far from perfect for providing comprehensive guarantees for all flows. Ideally, it is expected to monitor *every* flow with *minimum* errors, which empowers various fine-grained network management operations such as responsive diagnosis [17, 51, 67] and bounds for network localization [3, 50]. However, the bounds in existing algorithms are designed for specific traffic statistics such as heavy hitters or flow distributions. They are too coarse-grained to poor scalability or fails to reach the expected accuracy level.

CocoSketch
SIGCOMM '21

[source]

CocoSketch: High-Performance Sketch-based Measurement over Arbitrary Partial Key Query

Yinda Zhang^{1,3}, Zaoxing Liu¹, Ruixia Wang¹, Tong Yang¹, Jizhou Li¹, Ruijie Miao¹, Peng Liu¹, Ruwen Zhang¹, Junchen Jiang³
¹Peking University, ²Boston University, ³University of Chicago

Abstract

Sketch-based measurement has emerged as a promising alternative to the traditional sampling-based network measurement approaches due to its high accuracy and resource efficiency. While there have been various designs around sketches, they focus on measuring one particular flow key, and it is infeasible to support many keys based on these sketches. In this work, we take a significant step towards supporting *arbitrary partial key query*, where we only need to support a few partial keys for the whole network before measurement starts, and in query time, we can extract the information of any key in that range. We design CocoSketch, which casts arbitrary partial key queries to the sum subset estimation problem and makes the theoretical tools for subset sum estimation practical. To realize the desirable resource-accuracy tradeoffs in software-defined networking platforms, we propose two techniques: (1) stochastic variance minimization to more efficiently reuse per-packet update delay, and (2) removing modular dependencies in the per-packet update logic to make the implementation hardware-friendly. We implement CocoSketch on four popular platforms (CPU, Open vSwitch, F4, and FPGA) and show that compared to baselines that use traditional single-key sketches, CocoSketch improves average packet processing throughput by 27.2x and accuracy by 10.4x when addresses the design of NZE sketch, which is never studied.

CCS CONCEPTS

• Networks → Network monitoring. Network measurement;

Sketch; Arbitrary Partial Key Query; P4; FPPGA

ACM Reference Format:
Yinda Zhang, Zaoxing Liu, Ruixia Wang, Tong Yang, Jizhou Li, Ruijie Miao, Peng Liu, Ruwen Zhang, Junchen Jiang. 2021. CocoSketch: High-Performance Sketch-based Measurement over Arbitrary Partial Key Query. In *ACM SIGCOMM '21 Conference (SIGCOMM '21)*, August 25–28, 2021, Virtual Event, USA. ACM, New York, NY, USA, 16 pages. <https://doi.org/10.1145/3452296.3472892>

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyright © 2021, the author(s). This work is licensed by others under a Creative Commons license. The specific rights of each contributor are stated in the Author's Declaration and the license note for each document. This work may be reproduced, displayed, or distributed in whole or in part by other means without prior permission or express written permission of the copyright holders, provided the full text is not changed in any way, is not sold in any manner, and provides a link to the original source.

An ideal system for arbitrary partial key queries should meet three requirements: (1) **fidelity** (provable accuracy guarantee on any partial keys), (2) **resource efficiency** (minimizing the number of partial keys required to support an arbitrary partial key query), and (3) **compatibility** (on various software and hardware platforms, e.g., Open vswitch [59], PISA [49], and FPGA [51]).

Unfortunately, existing solutions that might support arbitrary partial key queries fall short on at least one requirement, as summarized in Table 1. R-HHH [39] reduces the overhead of updating

[†]Co-primary authors; Yikai Zhao, Kaicheng Yang, and Zirui Liu. Corresponding authors: Tong Yang (yangtong@mails.pku.edu.cn) and Yi Wang (wangy@ust.hk).

[‡]In this paper, per-flow level information means per-flow per-hop

Sketches have limitations. They **compute statistical summaries** and favor elements with **high frequency**.

$$\Pr \left[\hat{x}_i - x_i \geq \varepsilon \|\mathbf{x}\|_1 \right] \leq \delta$$

estimation error *relative to sum of all elements*

Sketches have limitations. They **compute statistical summaries** and favor elements with **high frequency**.

Let $\varepsilon = 0.01$, $\|x\|_1 = 10000$ ($\Rightarrow \varepsilon \cdot \|x\|_1 = 100$)

Assume two flows x_a , x_b ,

with $\|x_a\|_1 = 1000$, $\|x_b\|_1 = 50$



Sketches have limitations. They **compute statistical summaries** and favor elements with **high frequency**.

Let $\varepsilon = 0.01$, $\|x\|_1 = 10000$ ($\Rightarrow \varepsilon \cdot \|x\|_1 = 100$)

Assume two flows x_a , x_b ,

with $\|x_a\|_1 = 1000$, $\|x_b\|_1 = 50$

Error relative to **stream size**: 1%

Sketches have limitations. They **compute statistical summaries** and favor elements with **high frequency**.

Let $\varepsilon = 0.01$, $\|x\|_1 = 10000$ ($\Rightarrow \varepsilon \cdot \|x\|_1 = 100$)

Assume two flows x_a , x_b ,

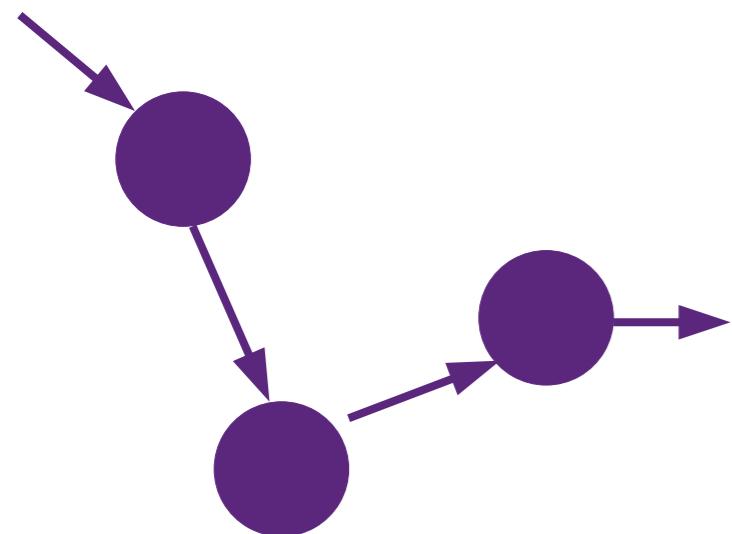
with $\|x_a\|_1 = 1000$, $\|x_b\|_1 = 50$

Error relative to **stream size**: 1%

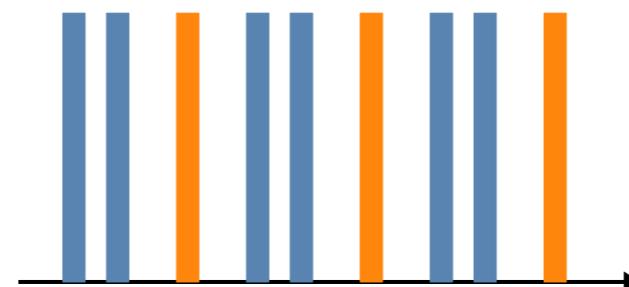
flow size: x_a : 10%, x_b : **200%**

Other problems a sketch can't handle:

causality



patterns



rare things



← You are looking at a stream of data (packets).
Today, I'll show you how set membership and
frequency queries can be realized in P4.

PART 1

Is a certain element (e.g. ip address) in the stream?
→ Bloom filter

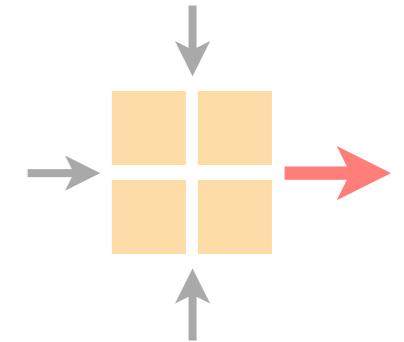
PART 2

How frequently does an element appear?
→ CountMin Sketch, Count Sketch, ...

TAKEAWAY

Probabilistic data structures provide
trade-offs between resources and error, and
provable guarantees to rely on.

Advanced Topics in Communication Networks



Laurent Vanbever
nsg.ee.ethz.ch

ETH Zürich
Tue 11 Oct 2022