

SNIA DEVELOPER CONFERENCE



BY Developers FOR Developers

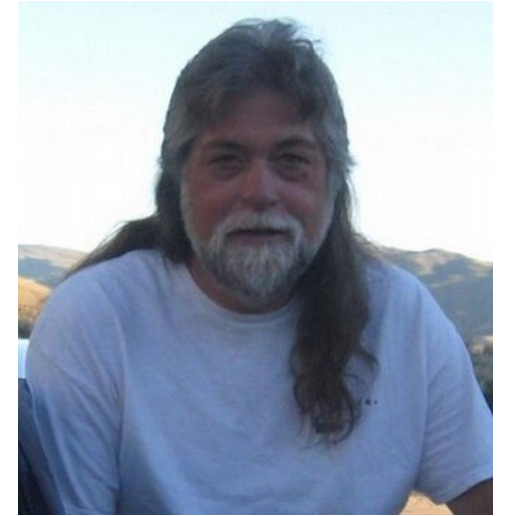
September 16-18, 2024
Santa Clara, CA

Exploring Optimizations of Content Delivery Networks

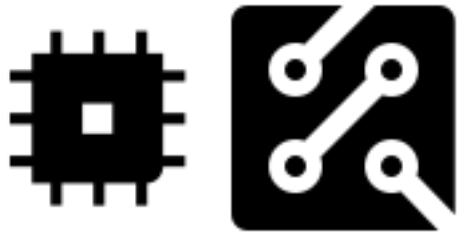
Andy Banta – Storage Janitor – Magnition IO

Andy Banta

Magnition.io (Consultant)
SolidFire (VMware development, acq. by NetApp)
DataGravity (Container exploitation lead)
VMware (iSCSI Tech Lead, IPO)
Sun Microsystems (Initial Fibre Channel development)
Patent, early distributed network projects, data acquisition
@andybanta

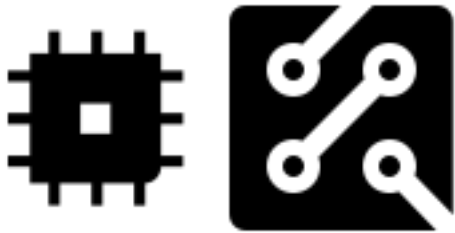


Engineering Simulation



Engineering Simulation

- Cheaper, faster, more flexible than system building
- Engineering design uses simulations, why not software?



How to Optimize All These Factors

- ✦ Compose simulations of complex memory and storage
- ✦ Break the simulation into components
- ✦ Allows the components to be assembled like building blocks
- ✦ Provide reasonable but constrained set of variables
- ✦ Run simulations with synthetic data or actual IO traces



Caching

- Hit-n-miss
- Promotion and demotion
- Complexity of tiering
 - How rapidly this becomes an unmanageable problem

Caching

- Hit or miss
 - Populate based on use or prediction
 - Variety of algorithms for lookup, allocation, eviction and aging
 - Can be tuned for workload

Cache

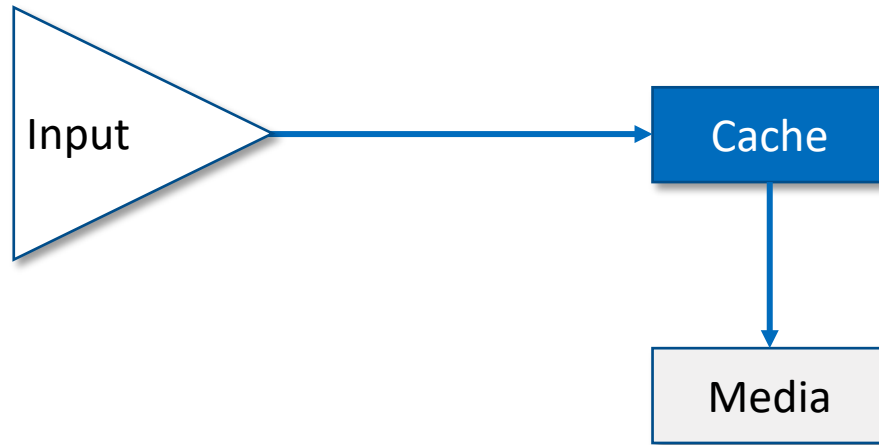
Caching

- Hit or miss
 - Populate based on use or prediction
 - Variety of algorithms for lookup, allocation, eviction and aging
 - Can be tuned for workload



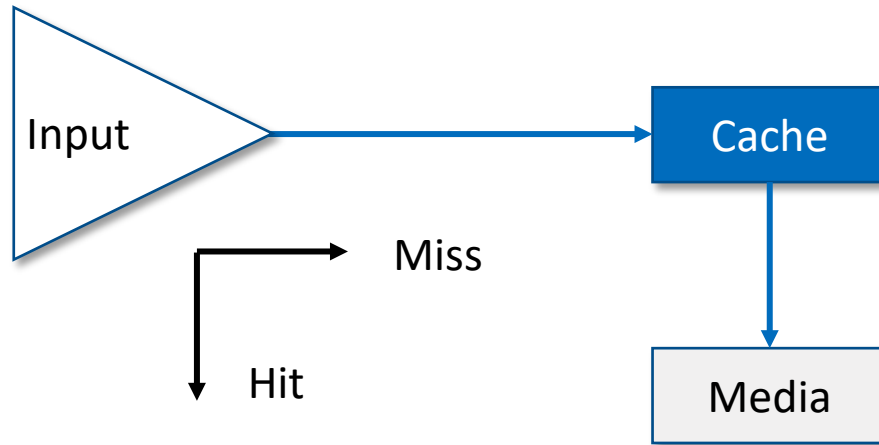
Caching

- Hit or miss
 - Populate based on use or prediction
 - Variety of algorithms for lookup, allocation, eviction and aging
 - Can be tuned for workload



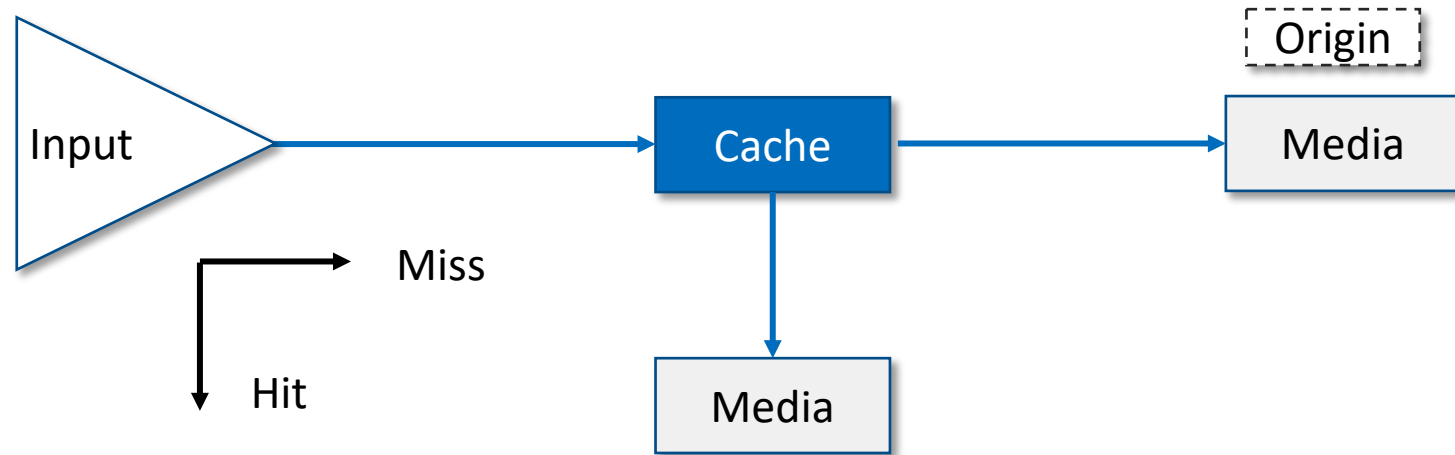
Caching

- Hit or miss
 - Populate based on use or prediction
 - Variety of algorithms for lookup, allocation, eviction and aging
 - Can be tuned for workload



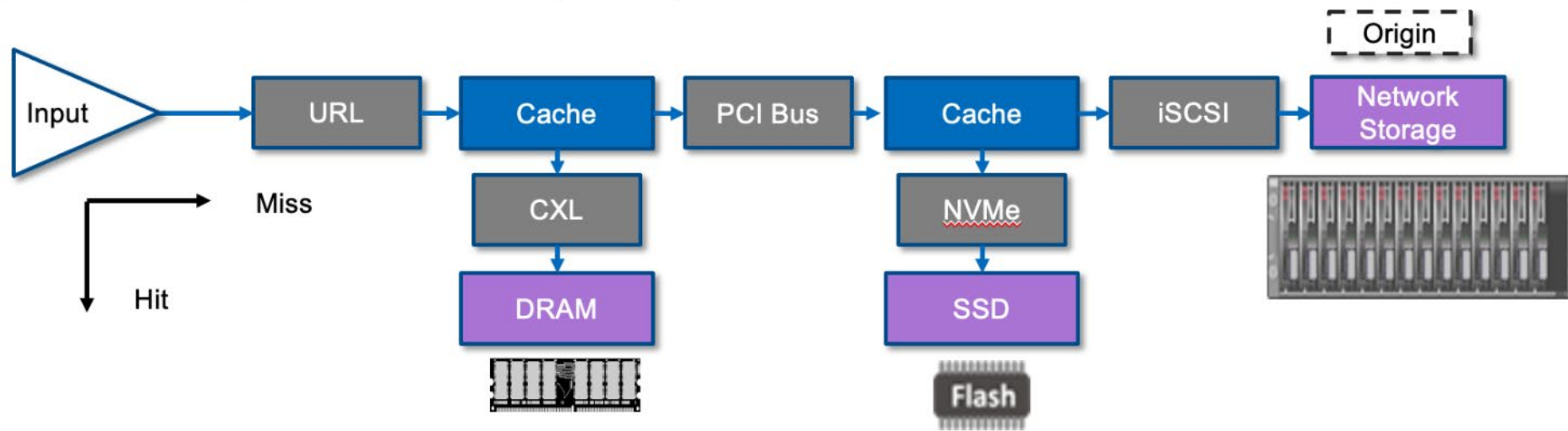
Caching

- Hit or miss
 - Populate based on use or prediction
 - Variety of algorithms for lookup, allocation, eviction and aging
 - Can be tuned for workload

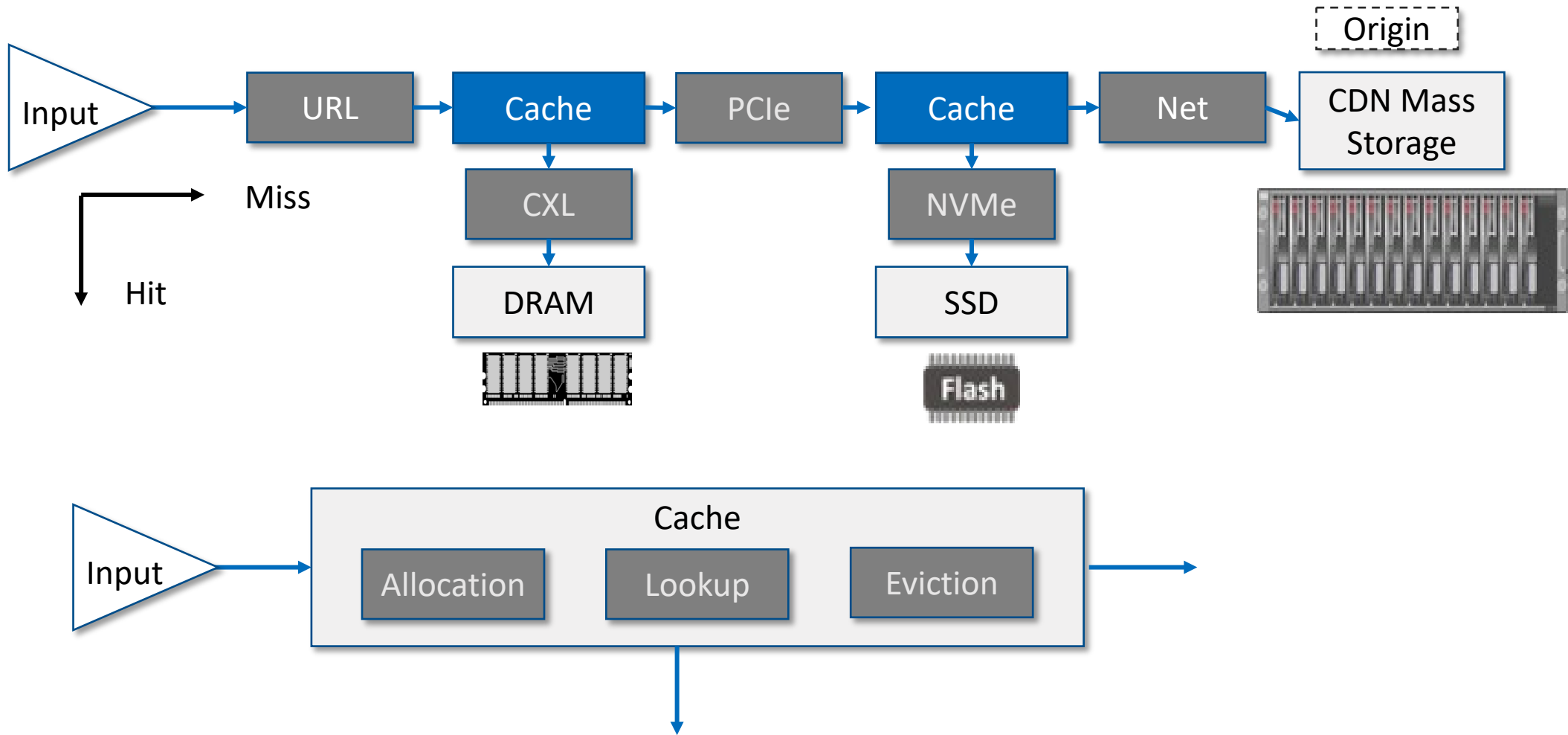


Multi-level Caching

- Currently used in Content Delivery Networks



Deterministic Behavioral Simulations



Time to Kick Out Old Eviction Algorithms

Introducing SIEVE

SIEVE: a New Eviction Algorithm

- Designed by Magnition intern
- Simple
- Efficient
- Improves scale and throughput

SIEVE operation

- Singly-linked list for eviction
- Pointer (“hand”) moves from tail to head, marks entries as unref’ed
- On cache hit, marks entry as referenced

An Example of SIEVE



SIEVE Operation

- On cache miss
 - Hand moves forward to next unreferenced entry and deletes it
 - New entry is added to the head of the list

An Example of SIEVE



SIEVE Operation

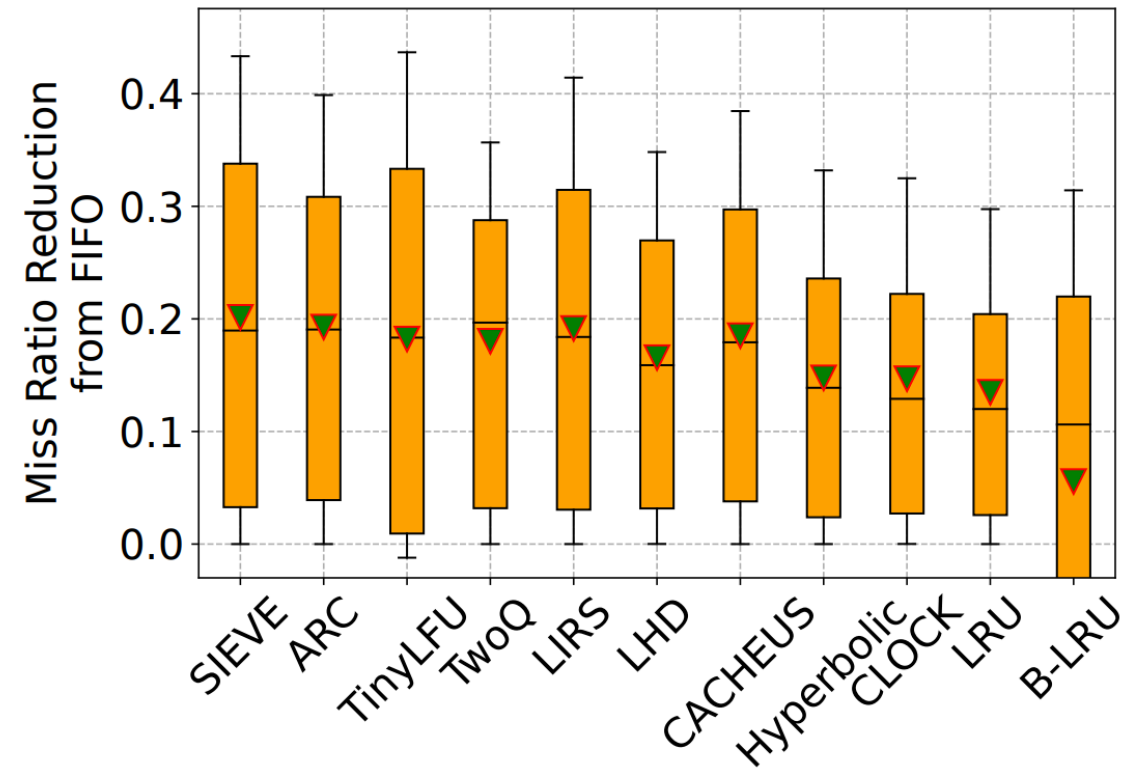
- If hand reaches head, it moves to tail

An Example of SIEVE



SIEVE Efficiency

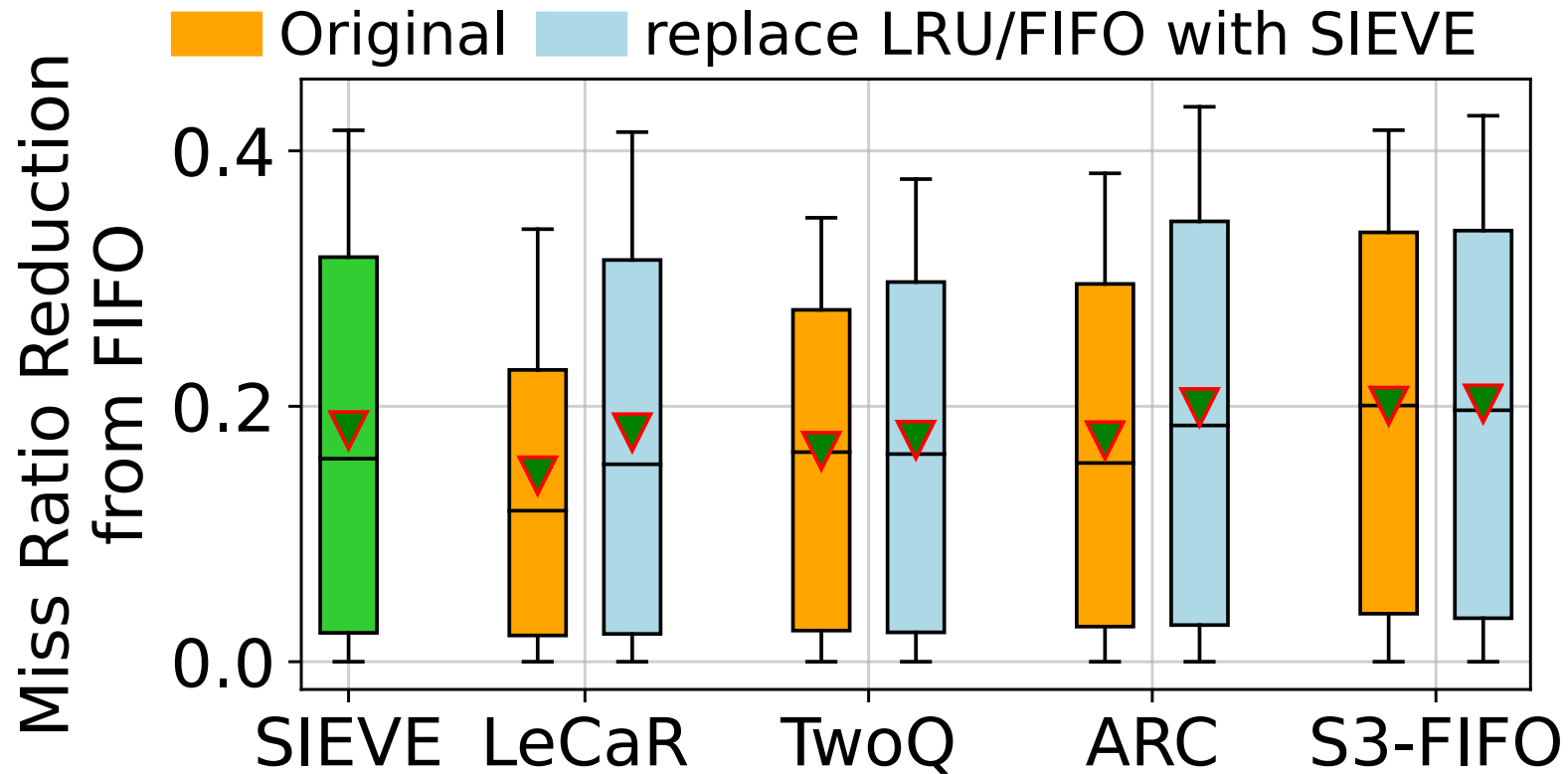
- Lower miss ratio in a variety of different applications



(a) CDN1 workloads, large cache, 1273 traces

SIEVE Efficiency

- SIEVE increases efficiency of more elaborate eviction algorithms



SIEVE Simplicity

- Modifications to existing algorithms are minimal
- Complete implementation is compact
 - Python sample implementation is 71 lines

Cache library	Language	Lines	Hour of work
groupcache	Golang	21	<1
mnemonist	Javascript	12	1
lru-rs	Rust	16	1
lru-dict	Python + C	21	<1

SIEVE Resources

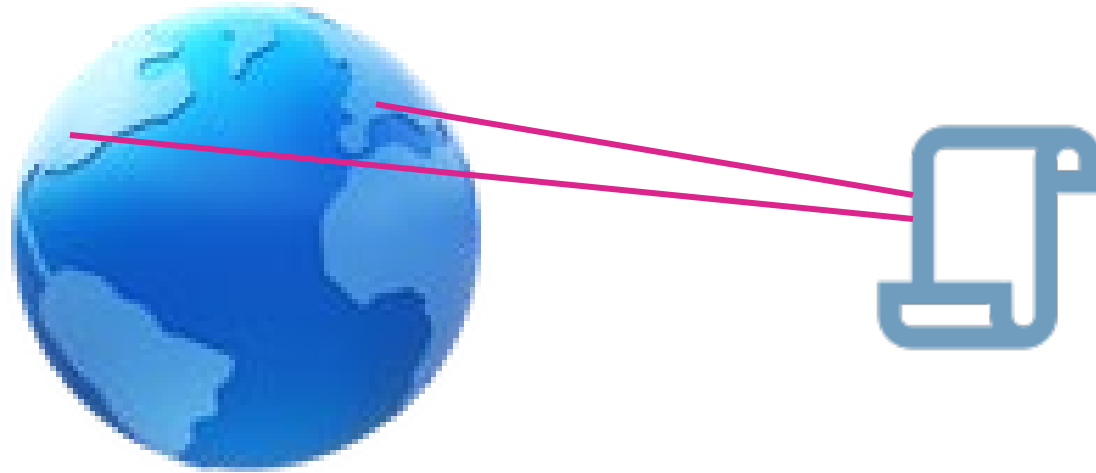
- Blog: <https://yazhuozhang.com/blog/2023/12/17/sieve-is-simpler-than-lru/>
- Github source: <https://cachemon.github.io/SIEVE-website/>
- Research paper: <https://junchengyang.com/publication/nsdi24-SIEVE.pdf>
- Shortcomings
 - Not scan resistant
 - Therefore not a good candidate for storage caches

Optimizing CDNs at Scale

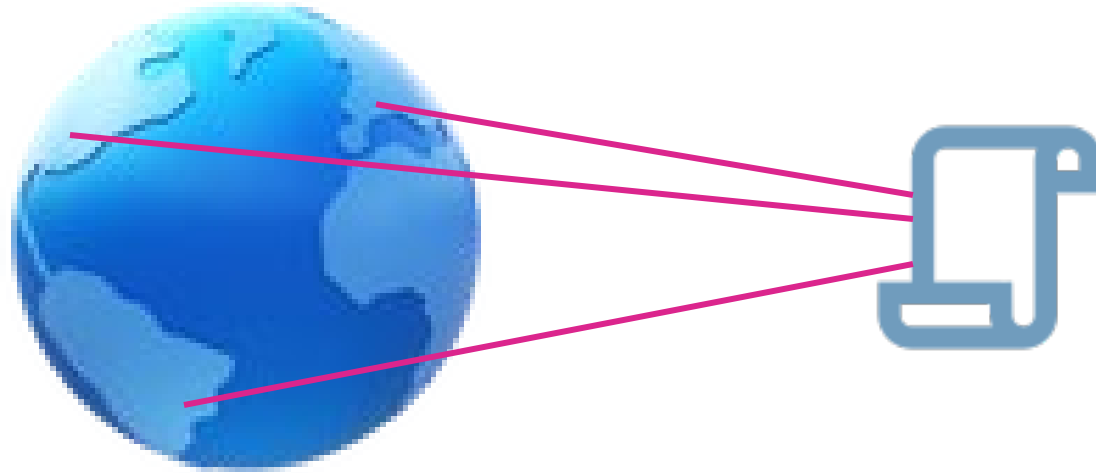
Content Delivery Networks



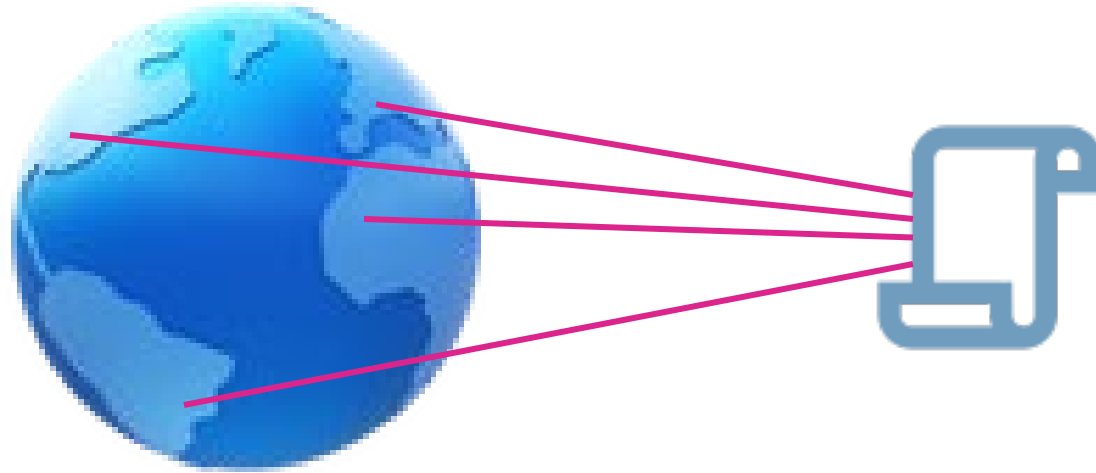
Content Delivery Networks



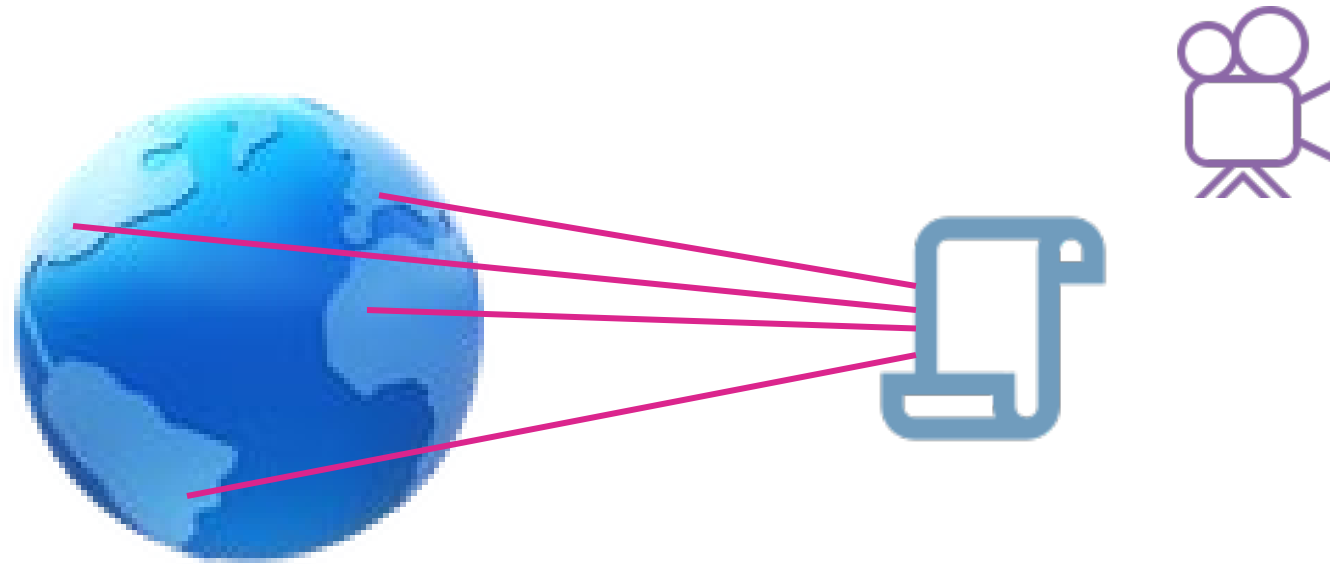
Content Delivery Networks



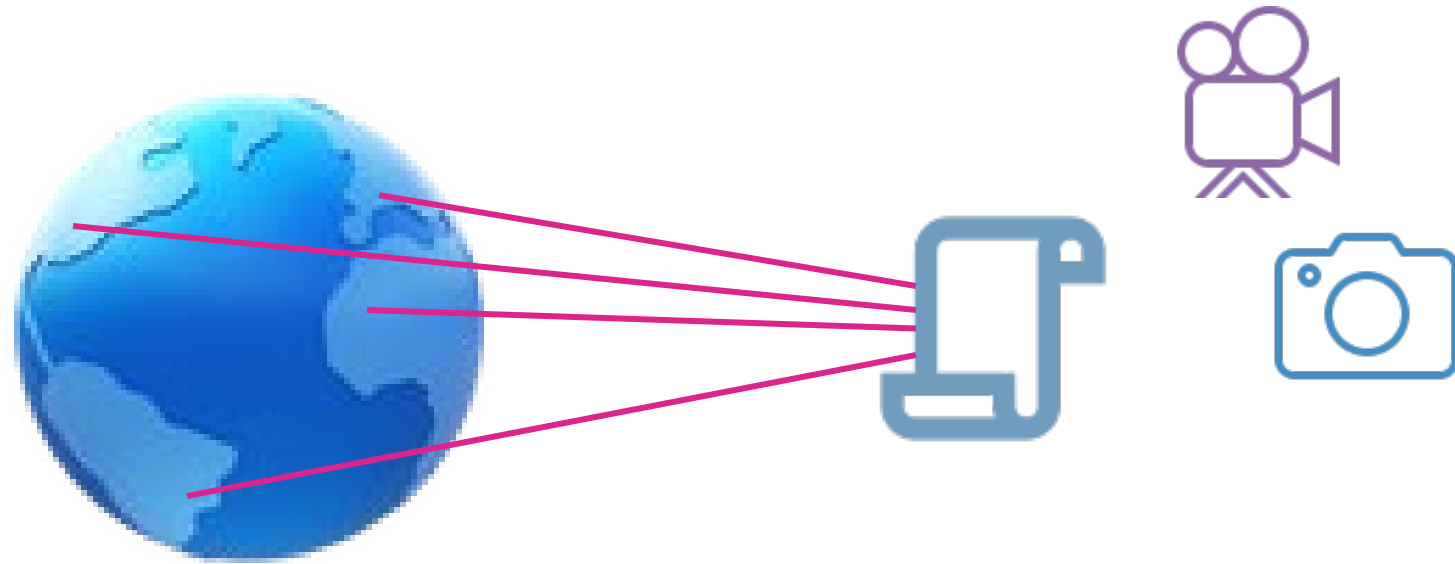
Content Delivery Networks



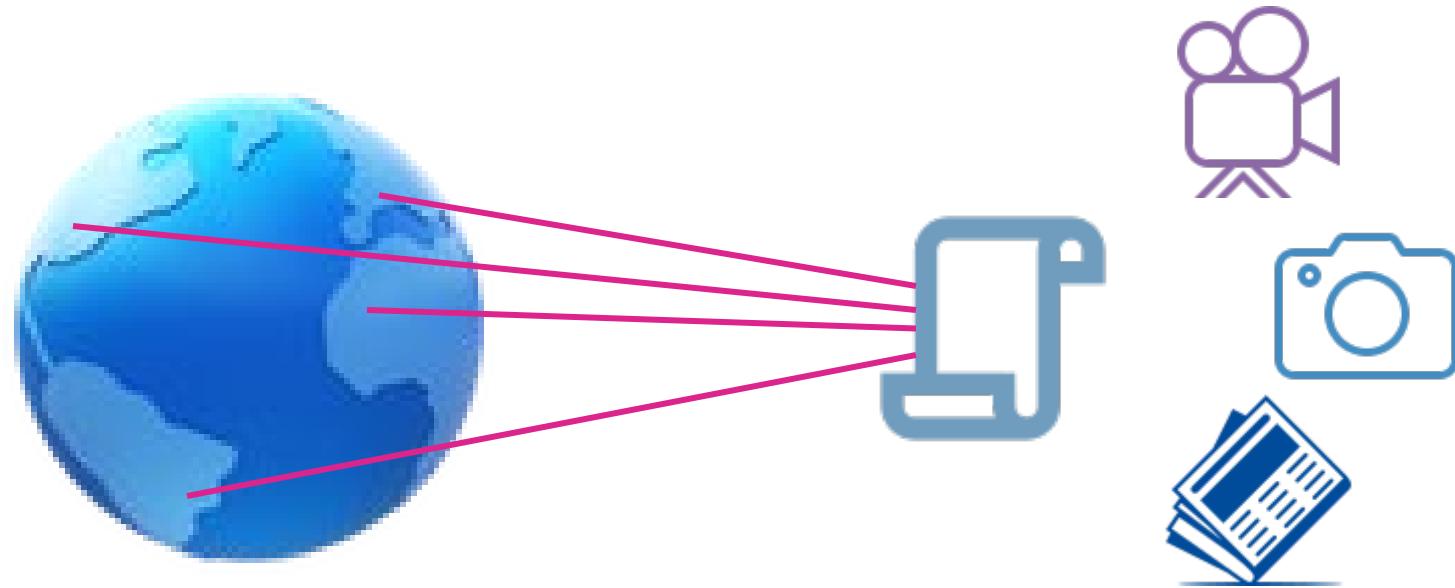
Content Delivery Networks



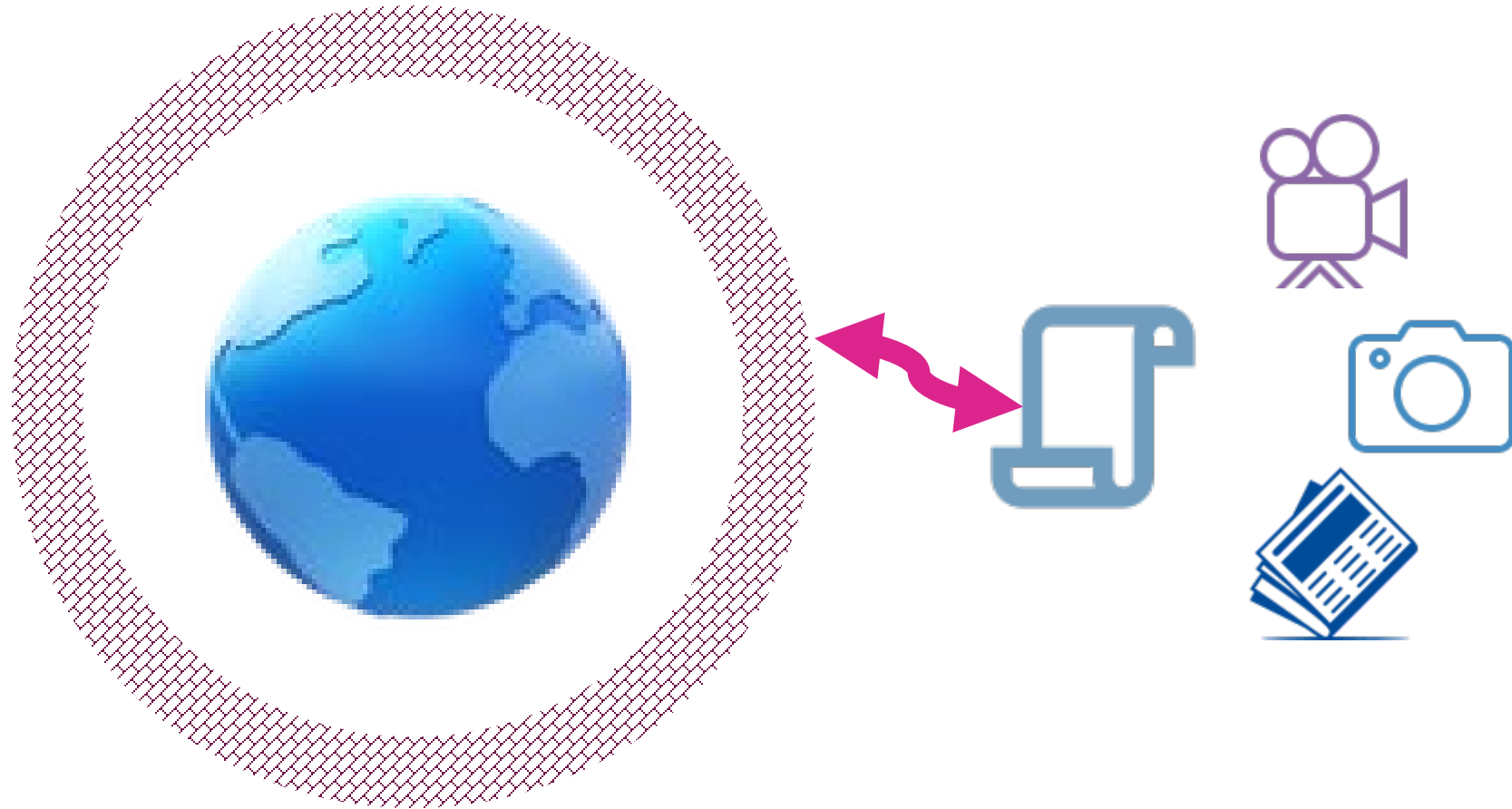
Content Delivery Networks



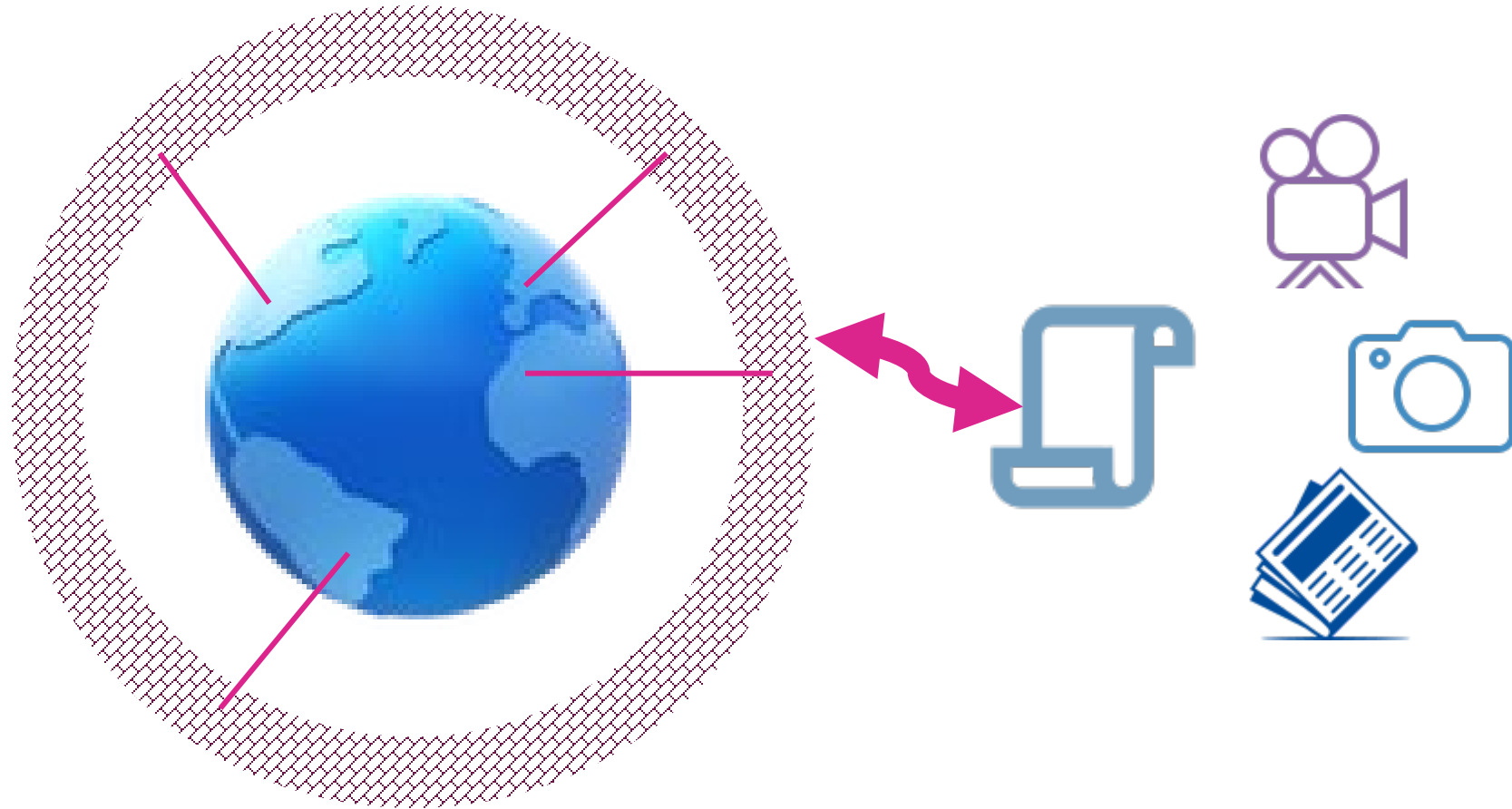
Content Delivery Networks



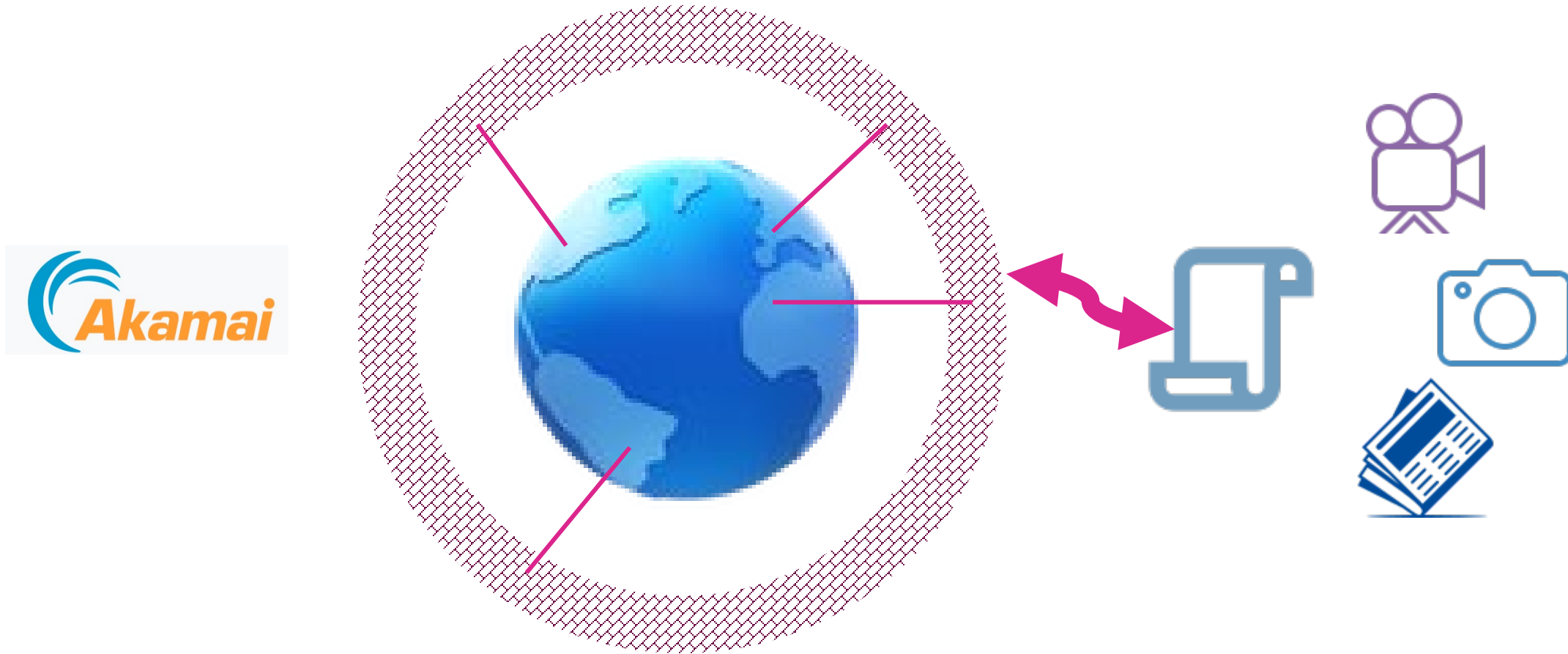
Content Delivery Networks



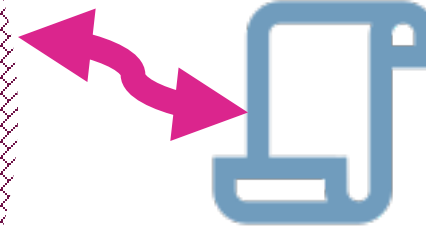
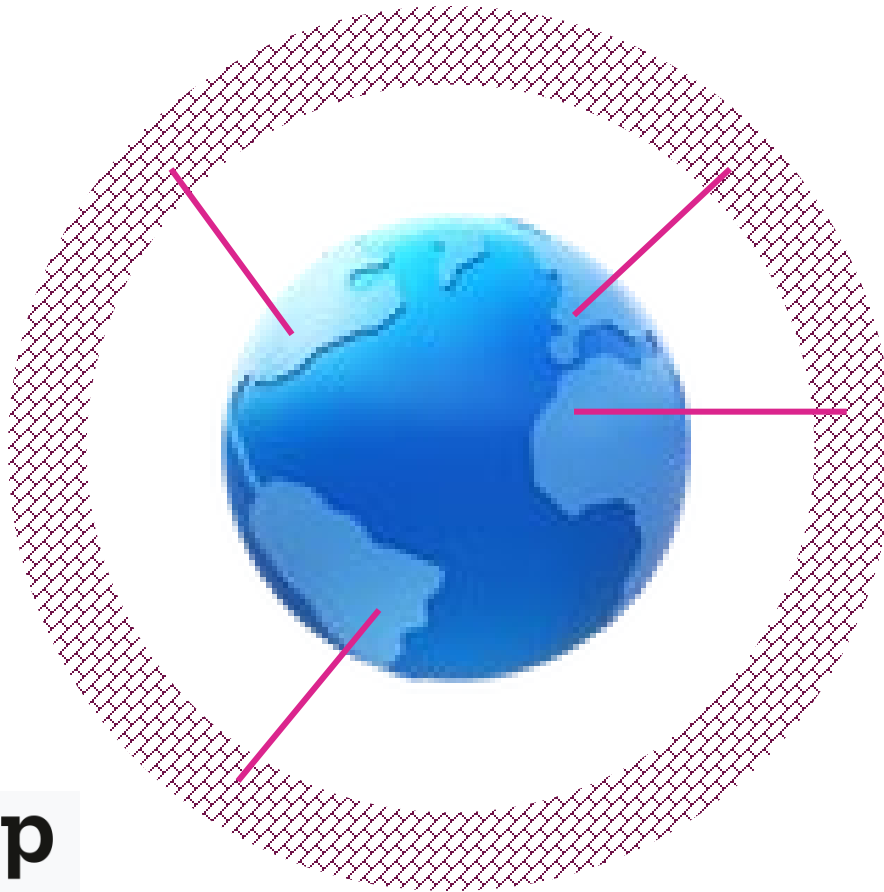
Content Delivery Networks



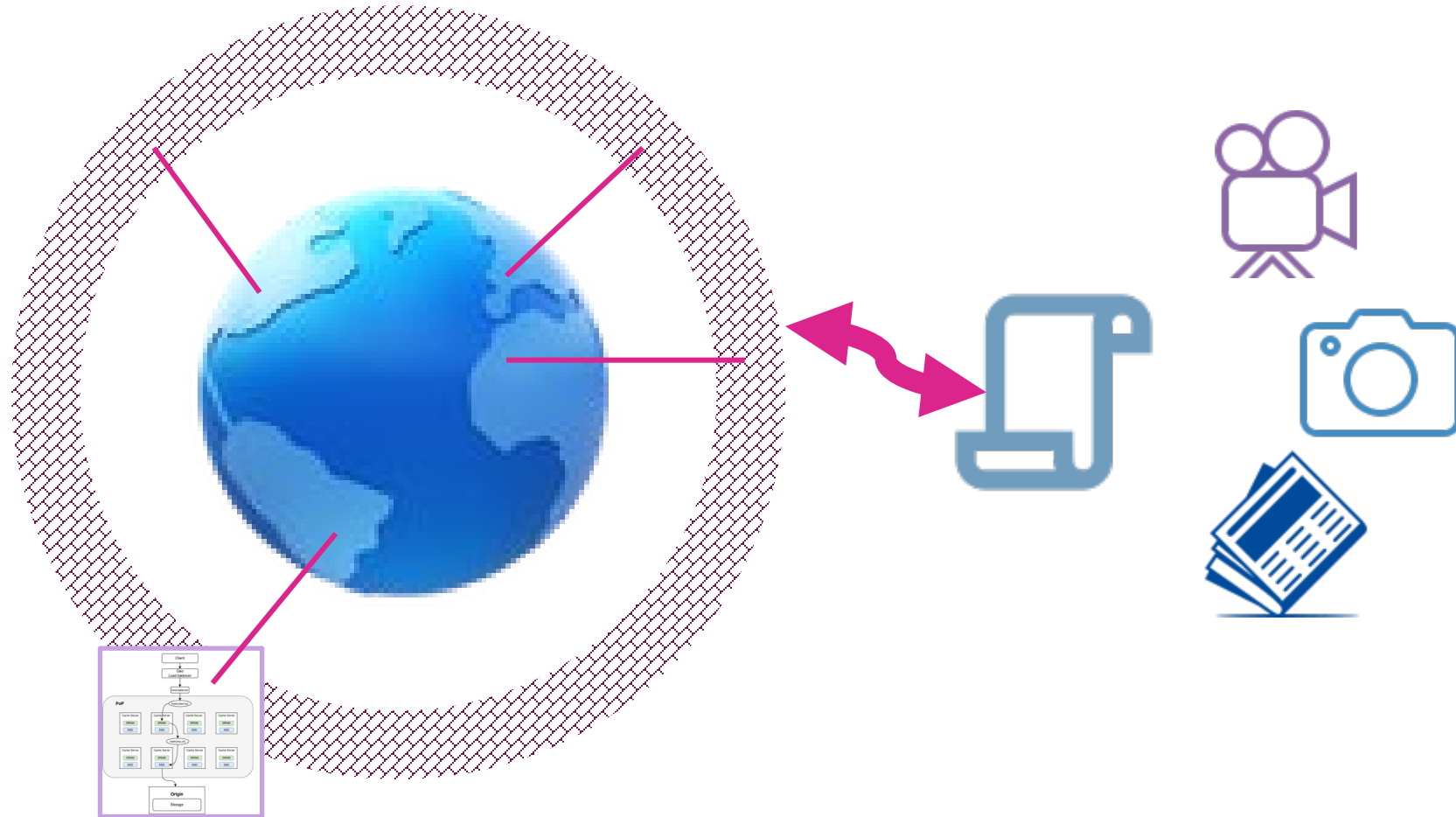
Content Delivery Networks



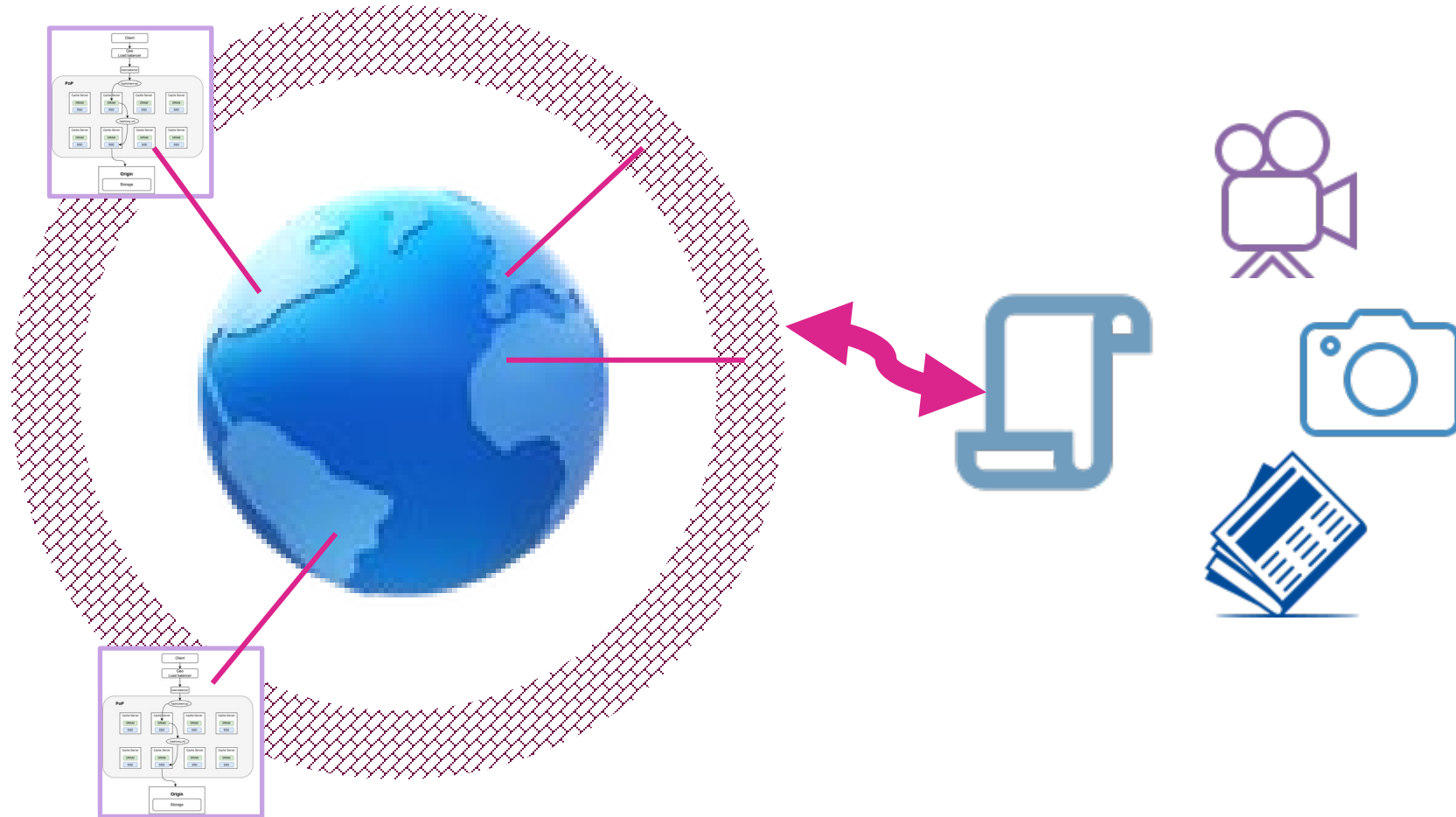
Content Delivery Networks



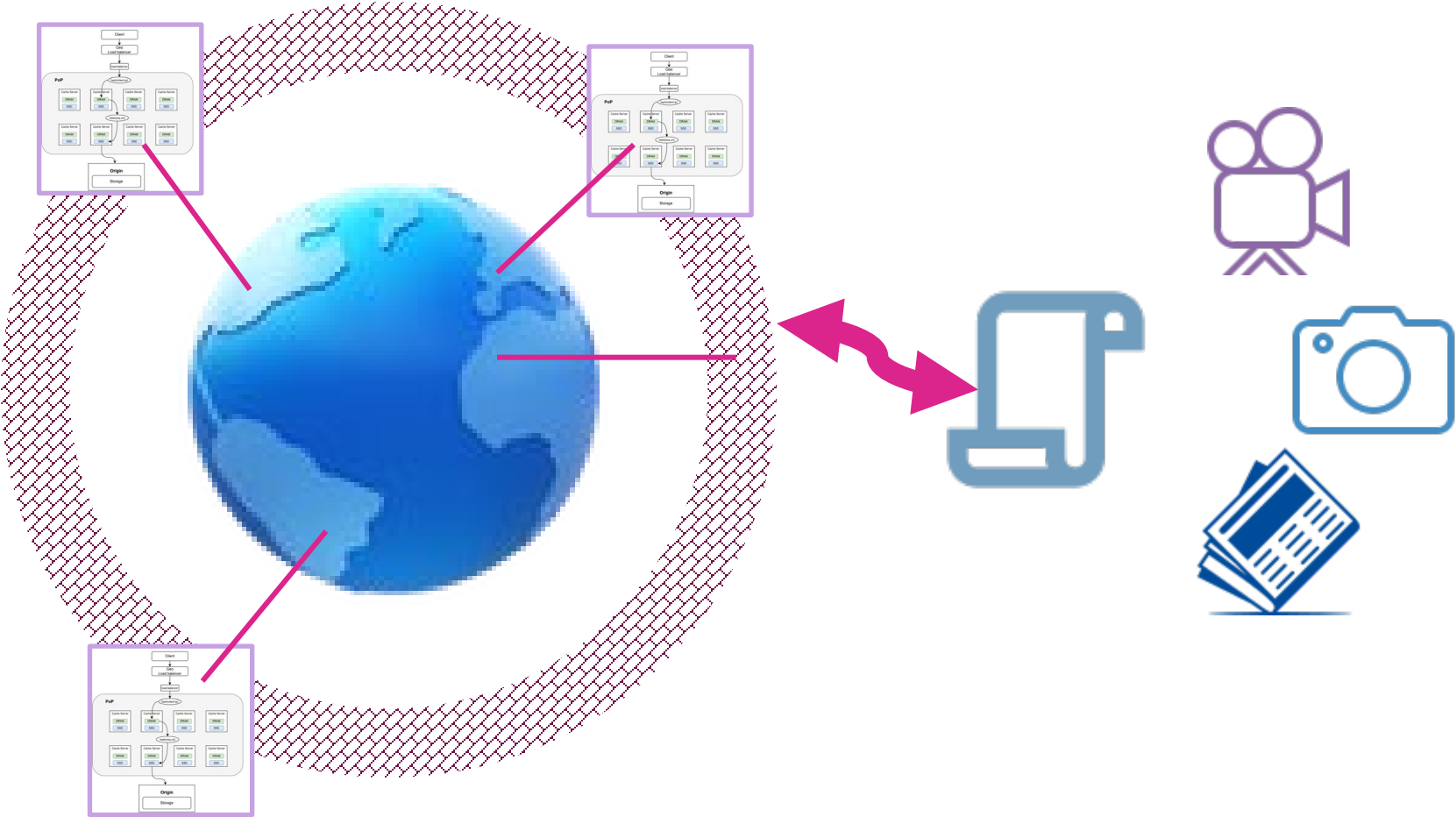
Content Delivery Networks



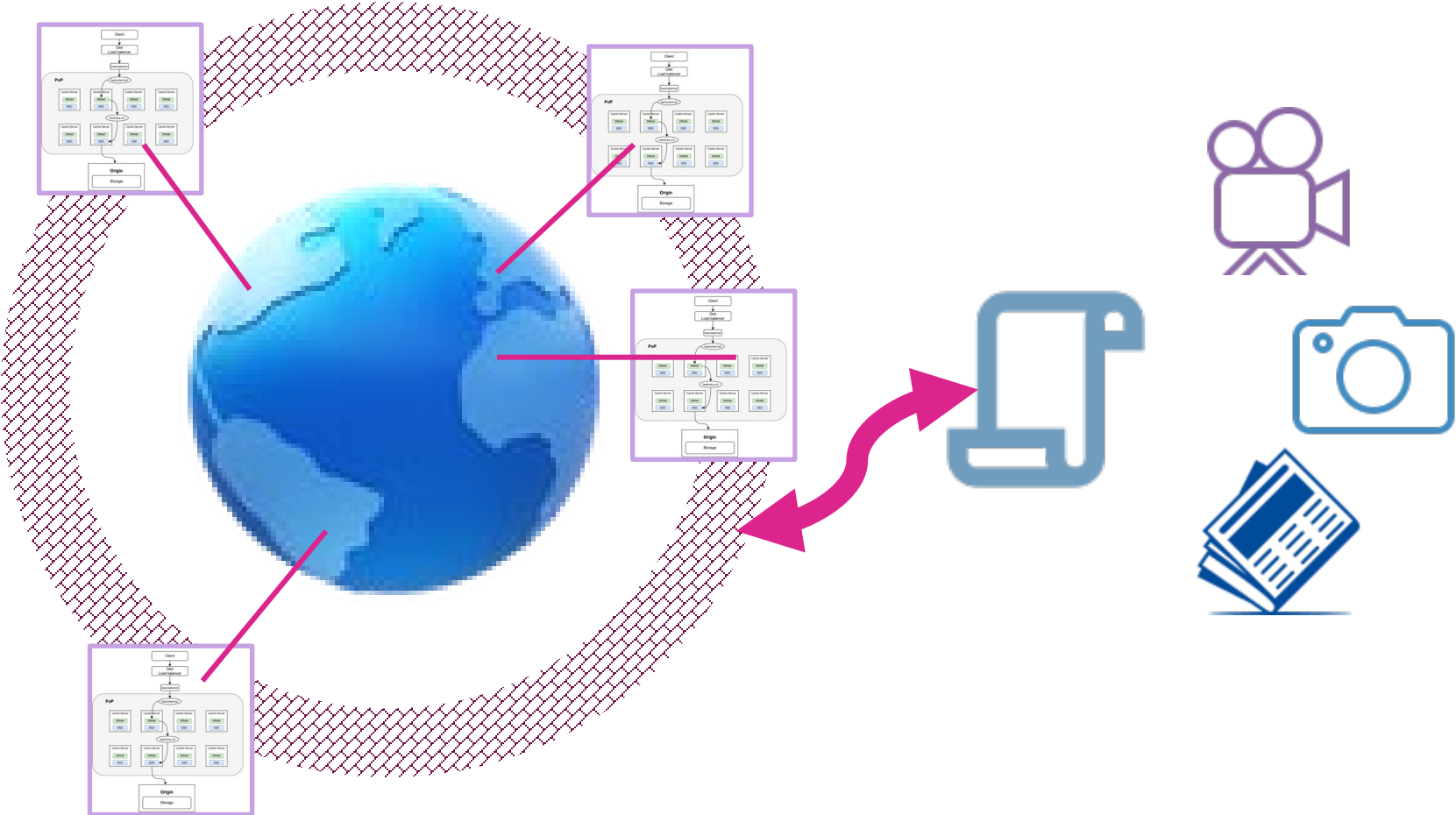
Content Delivery Networks



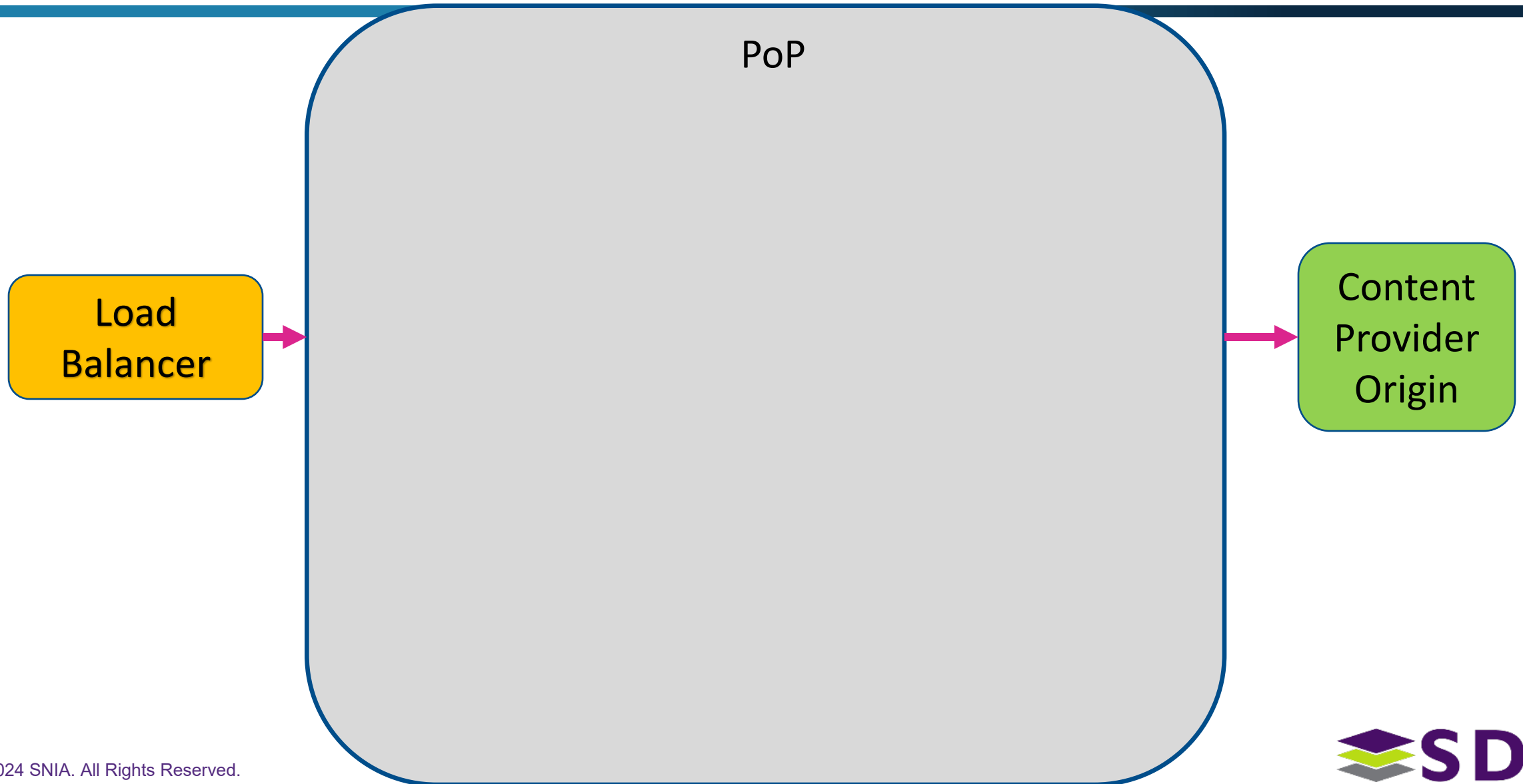
Content Delivery Networks



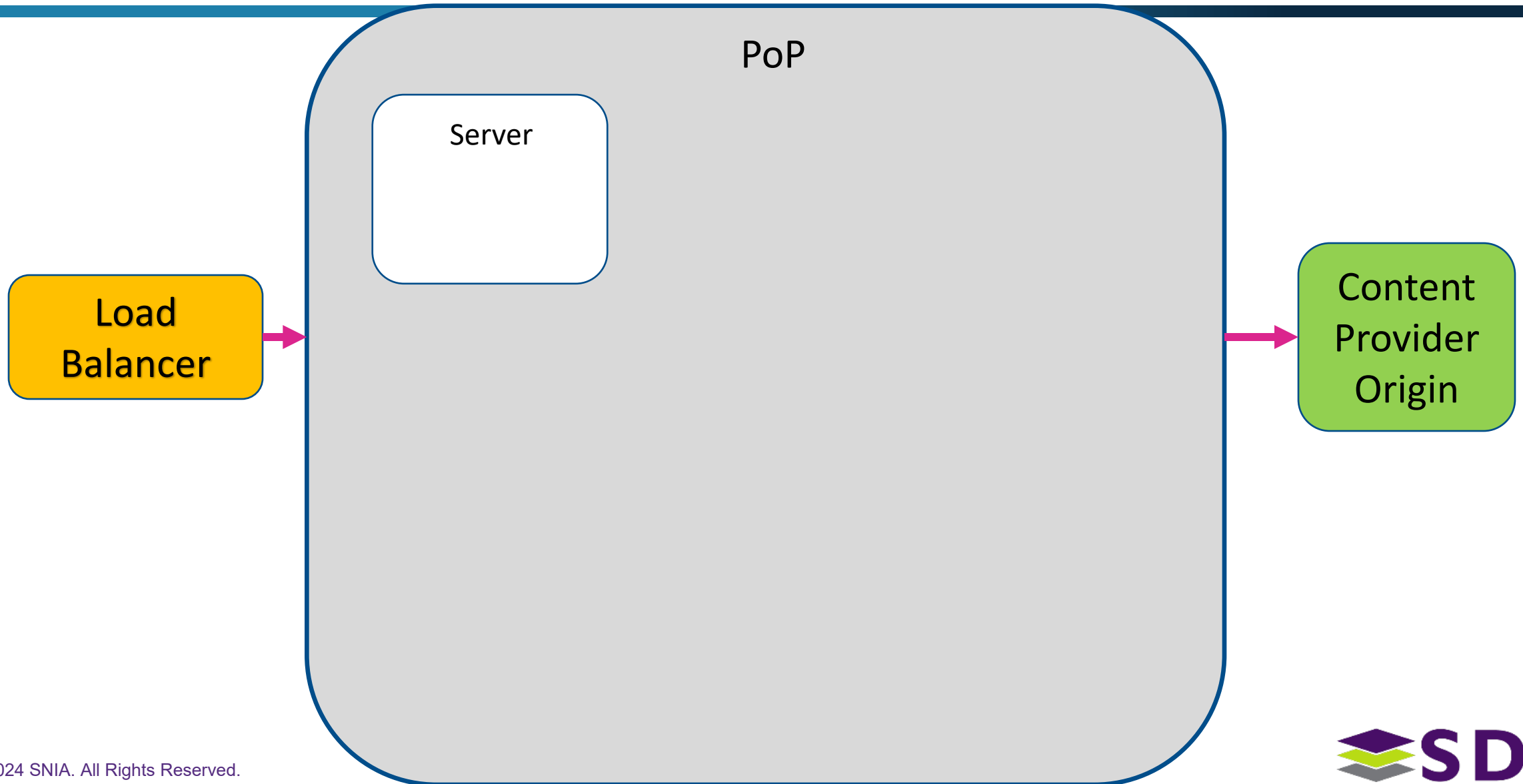
Content Delivery Networks



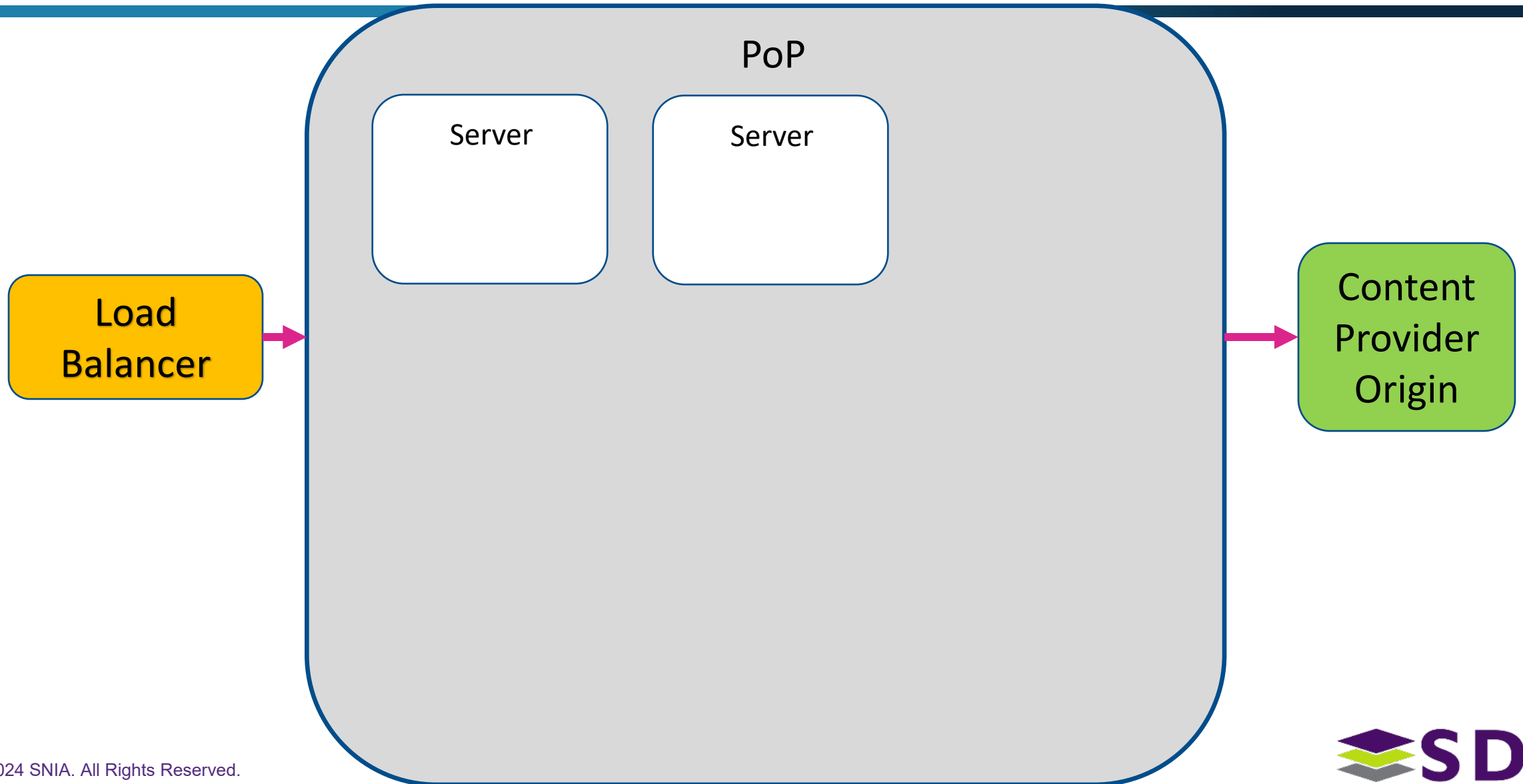
CDN PoP



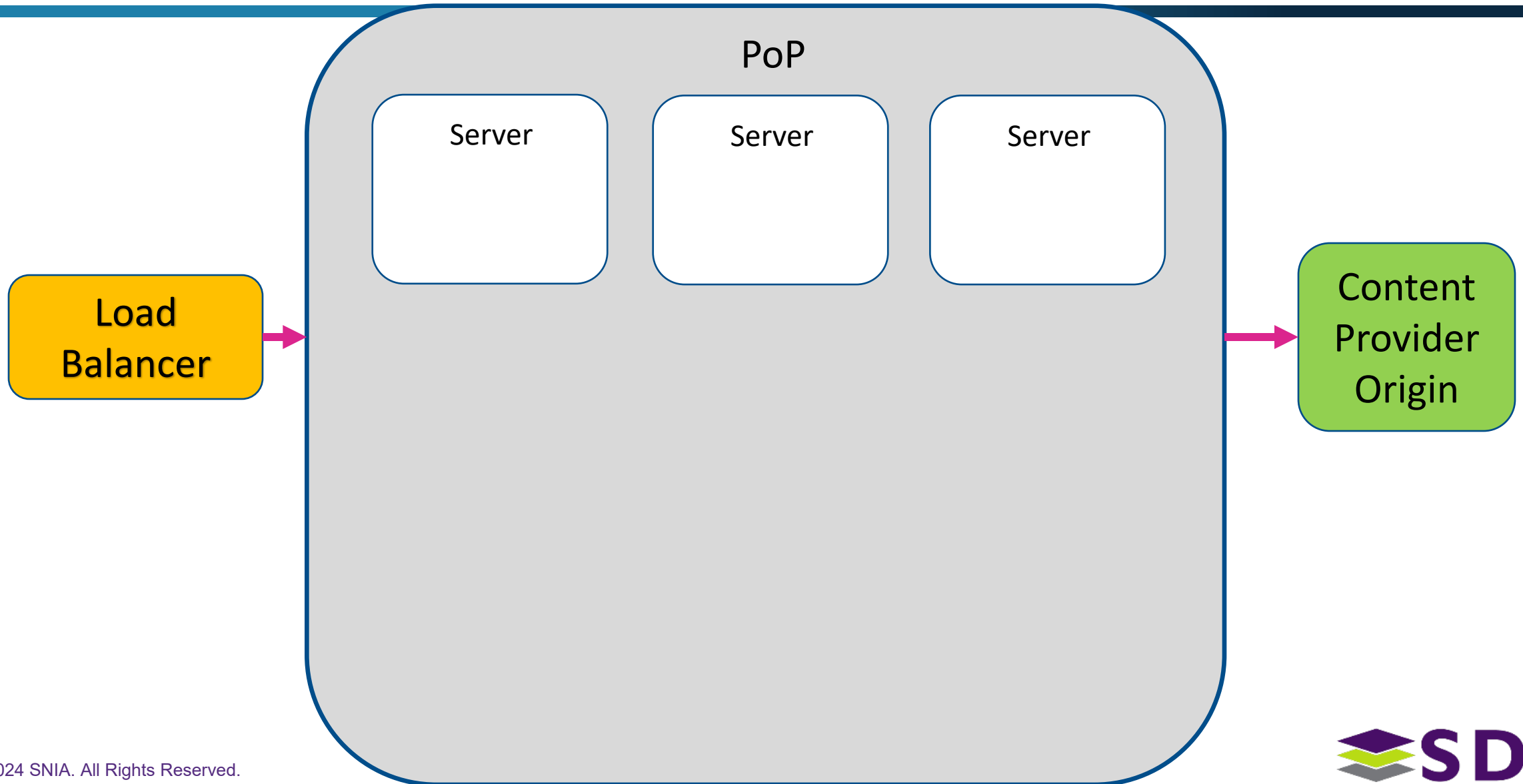
CDN PoP



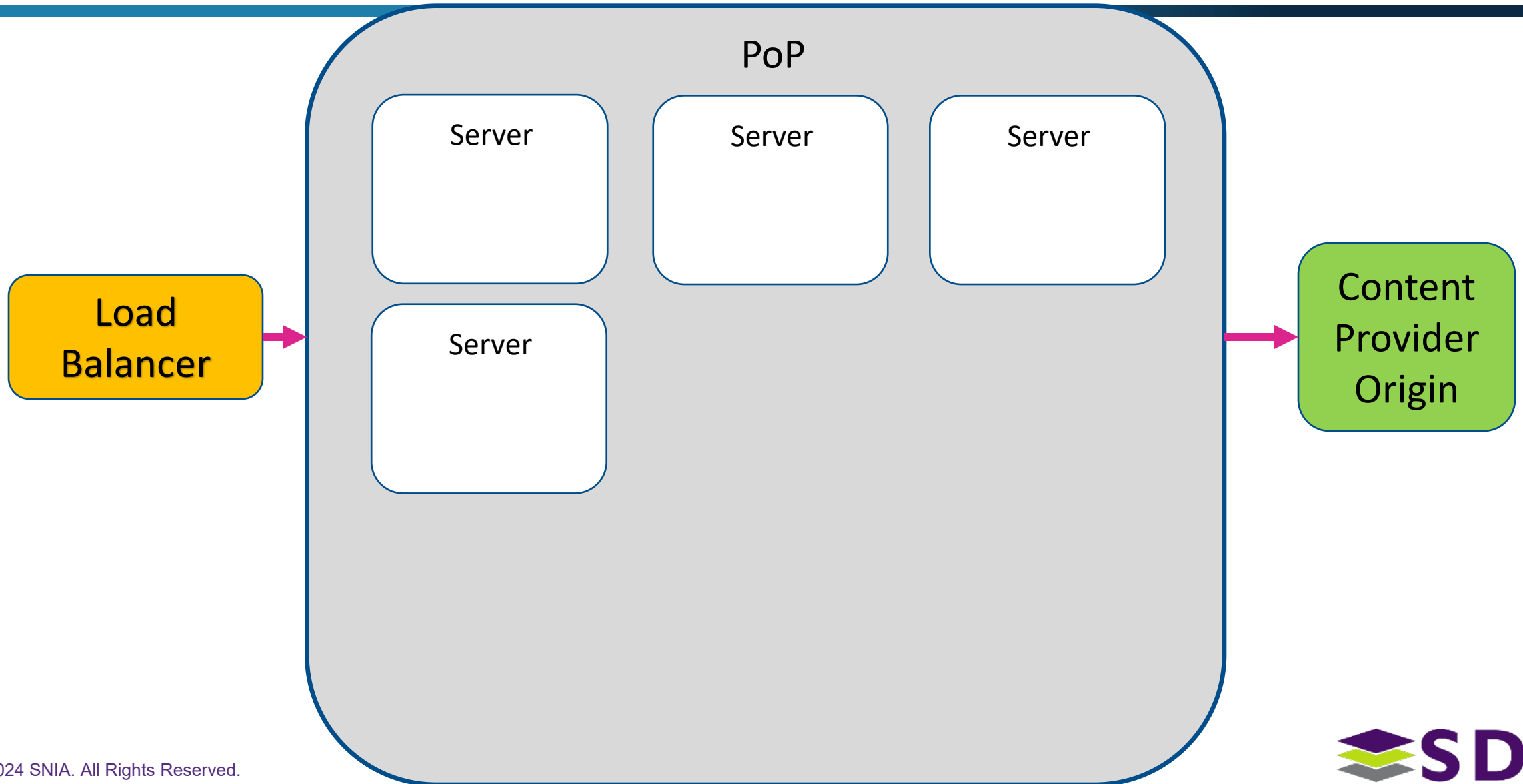
CDN PoP -- Compute



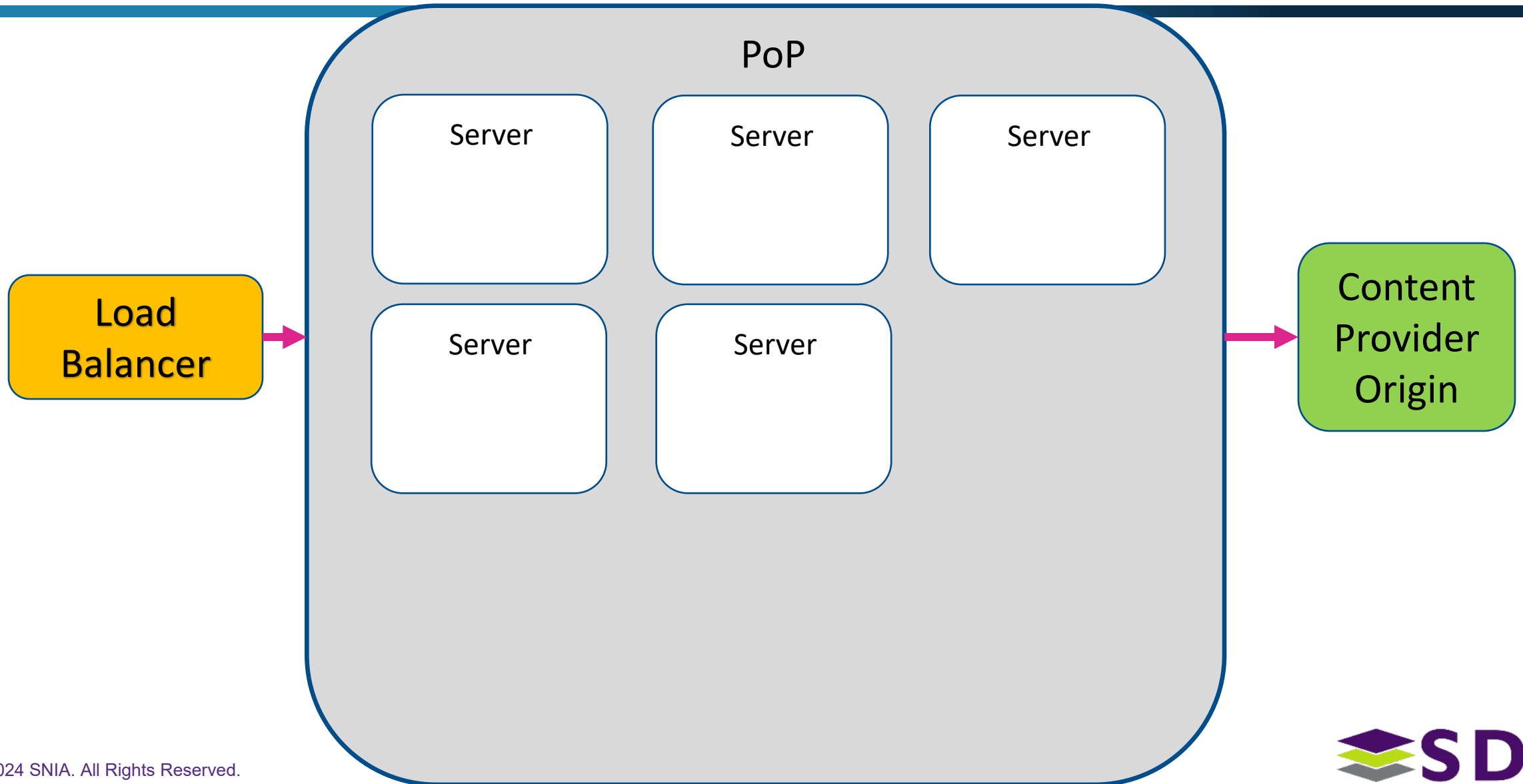
CDN PoP -- Compute



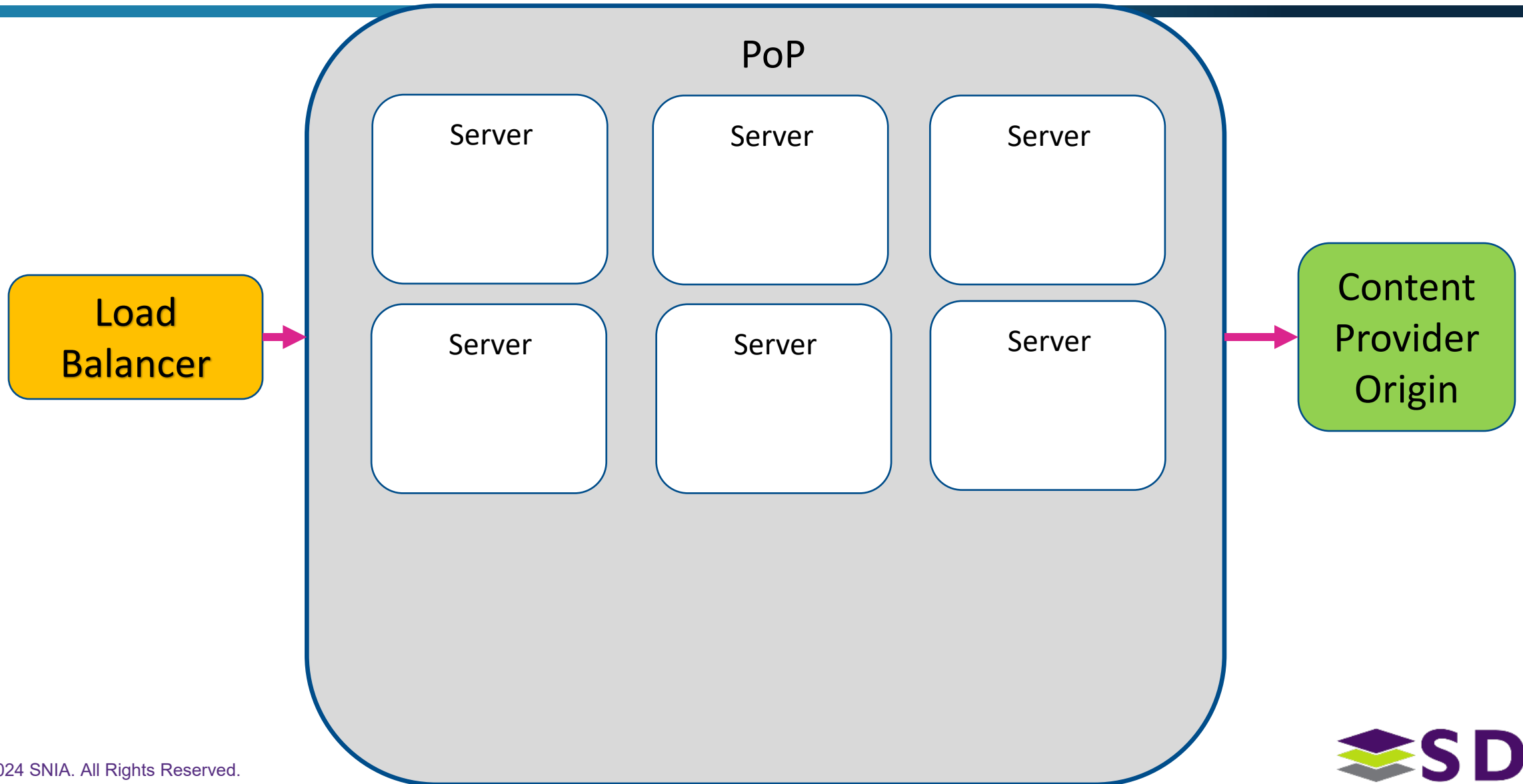
CDN PoP -- Compute



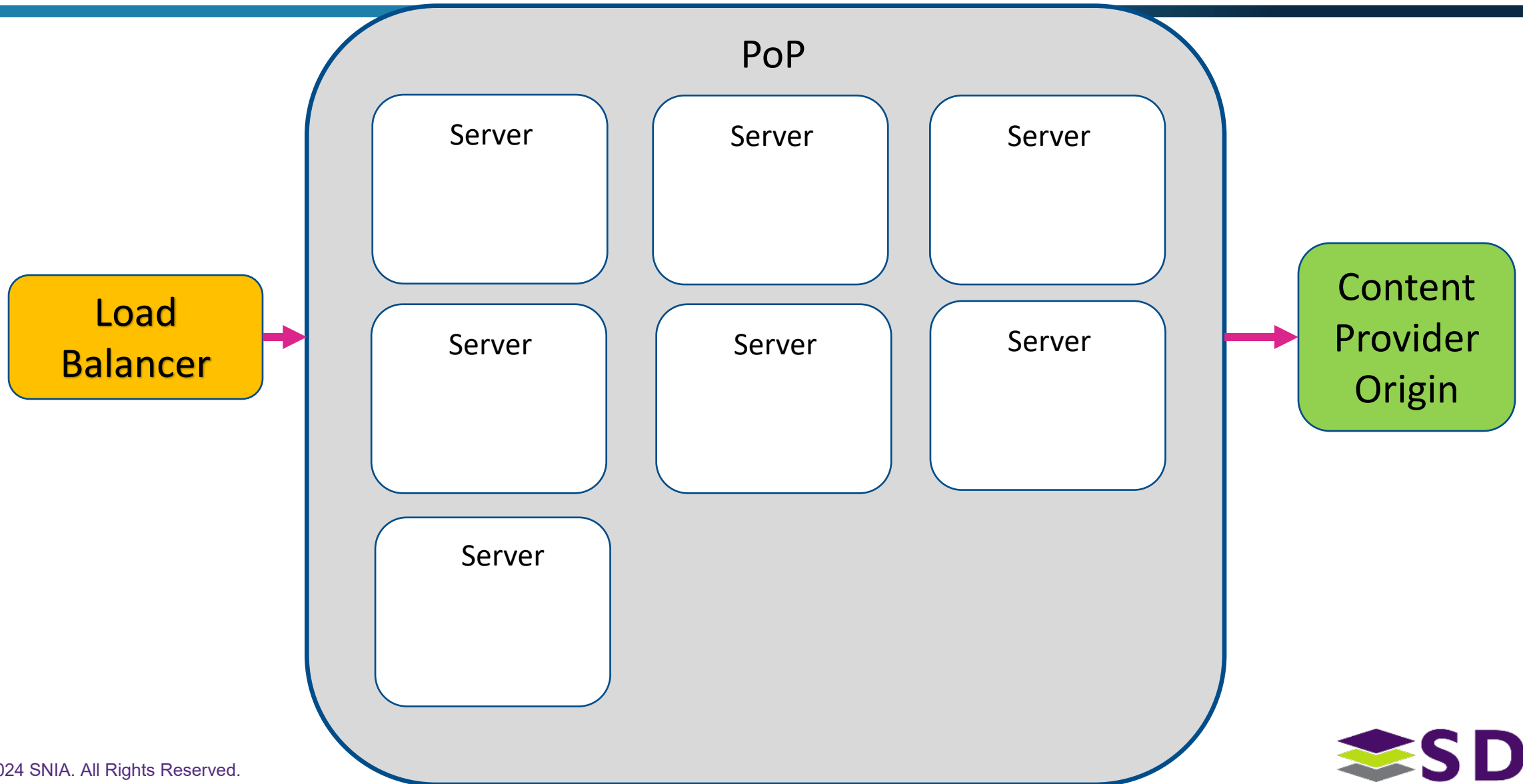
CDN PoP -- Compute



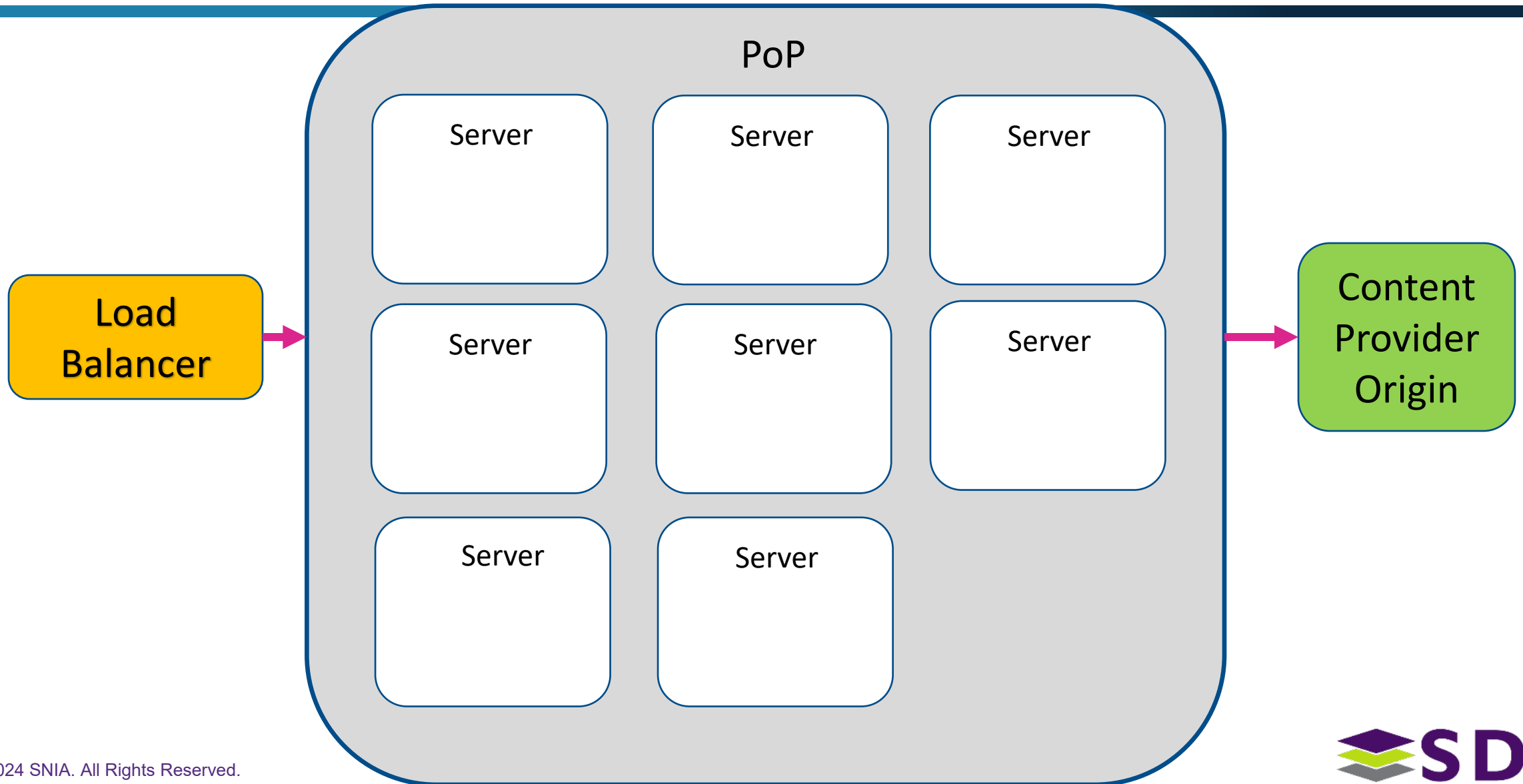
CDN PoP -- Compute



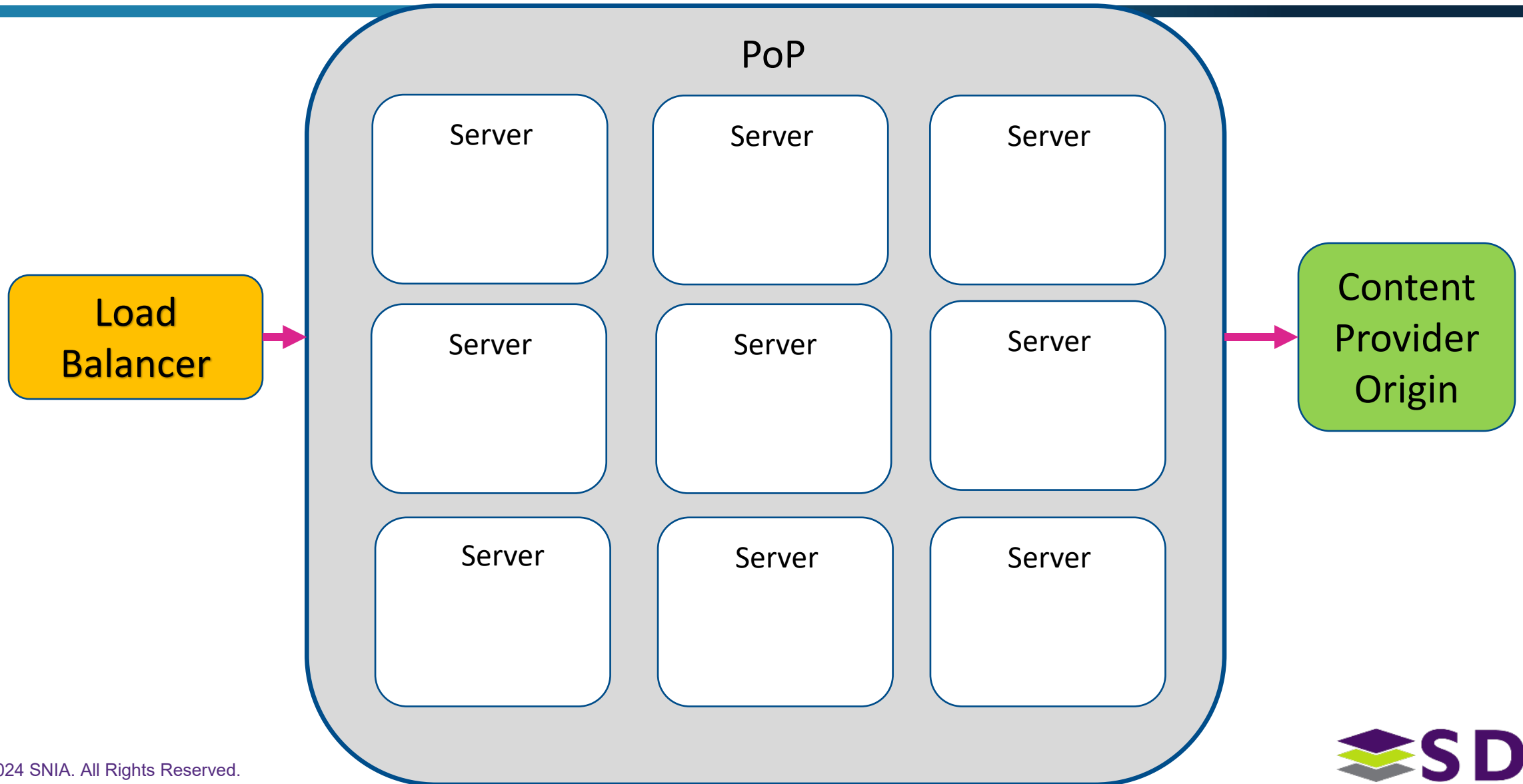
CDN PoP -- Compute



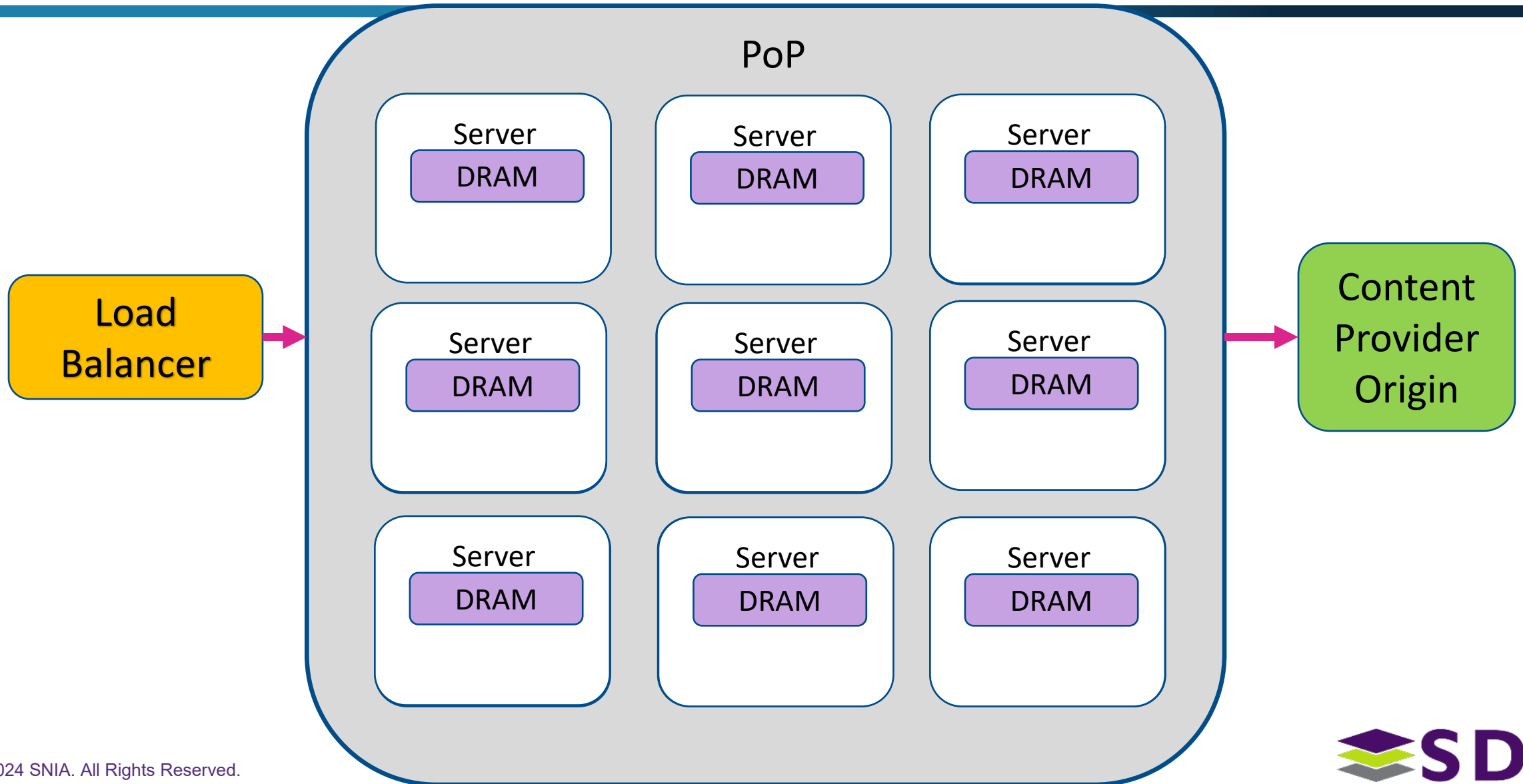
CDN PoP -- Compute



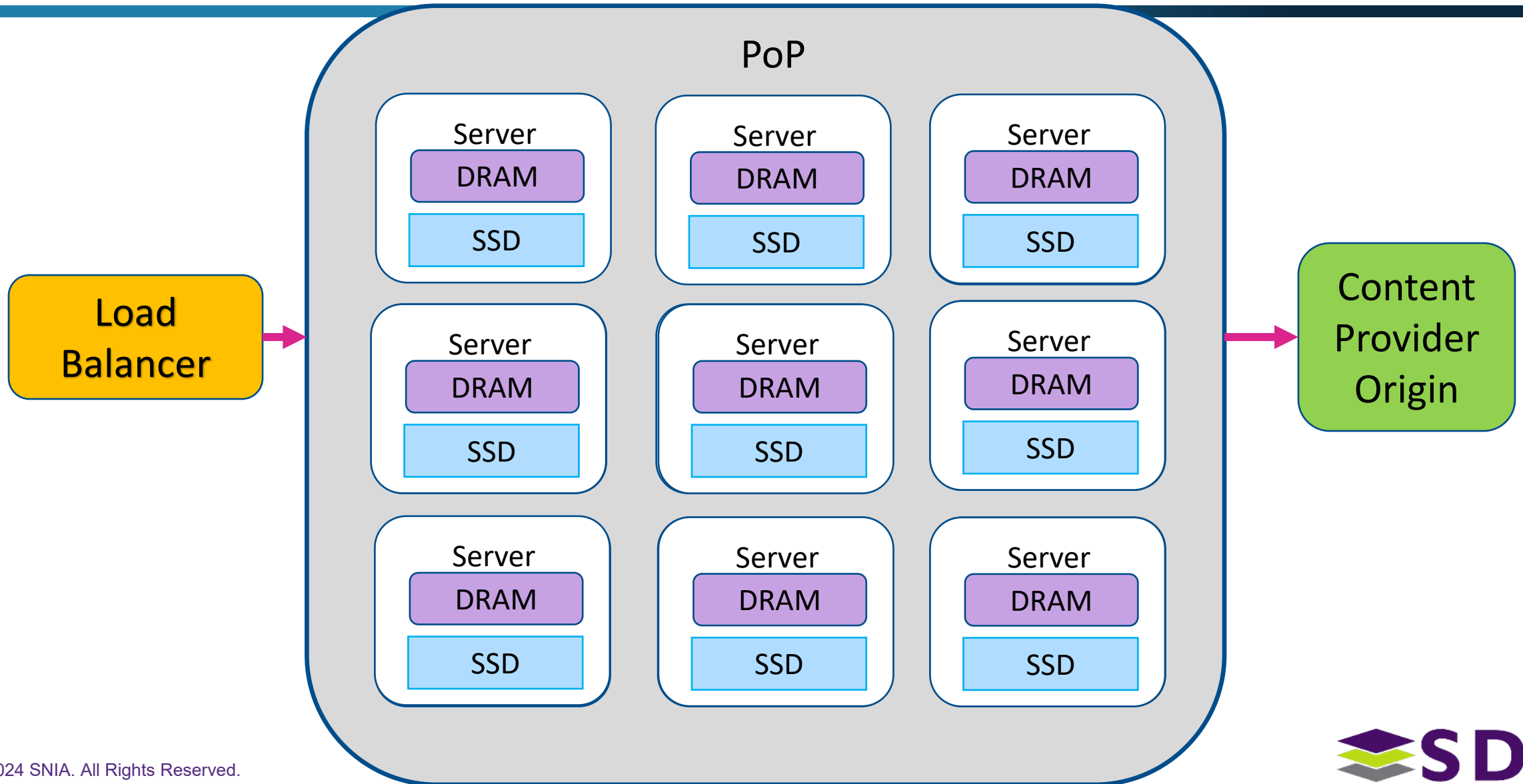
CDN PoP -- Compute



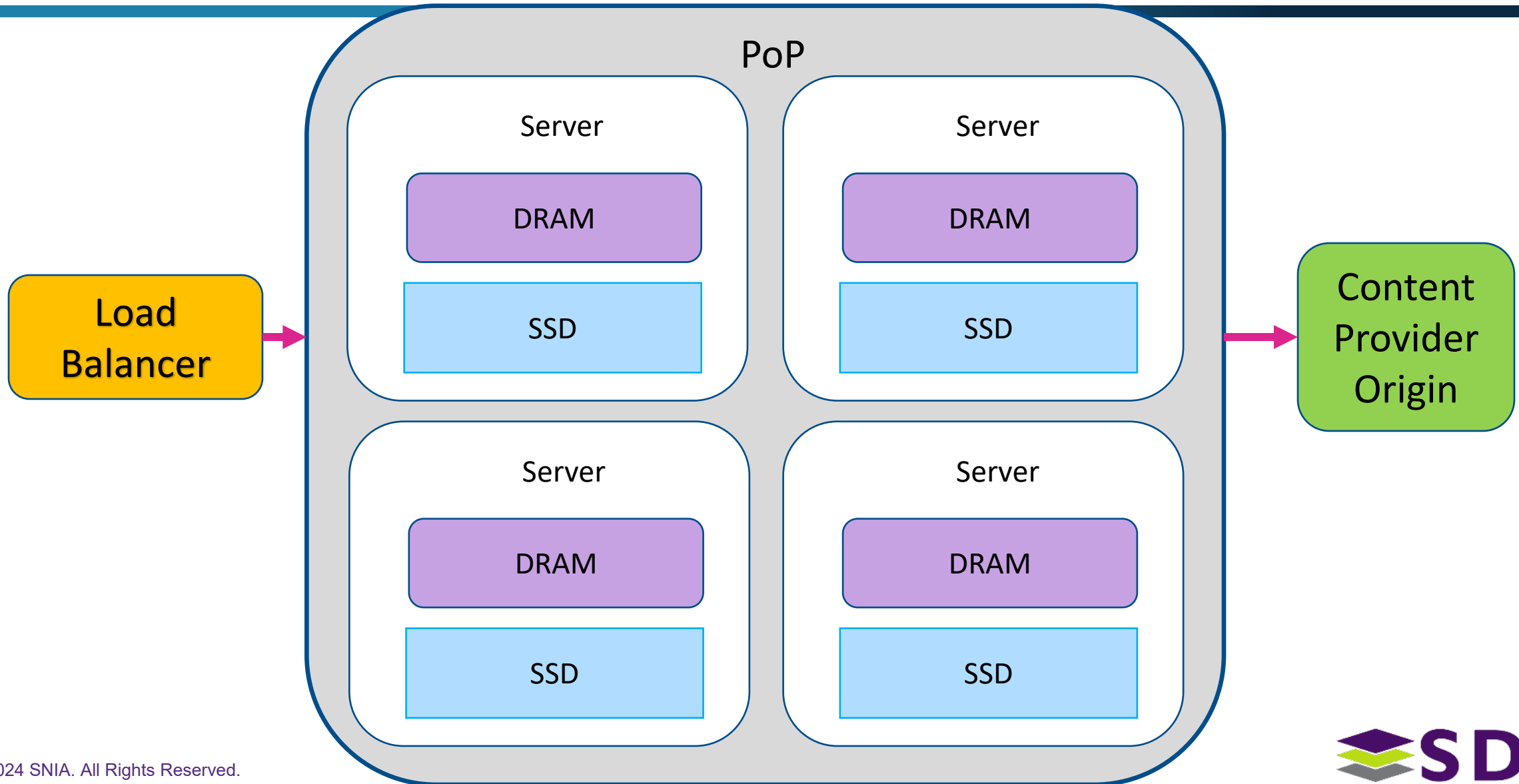
CDN PoP – Compute + Memory



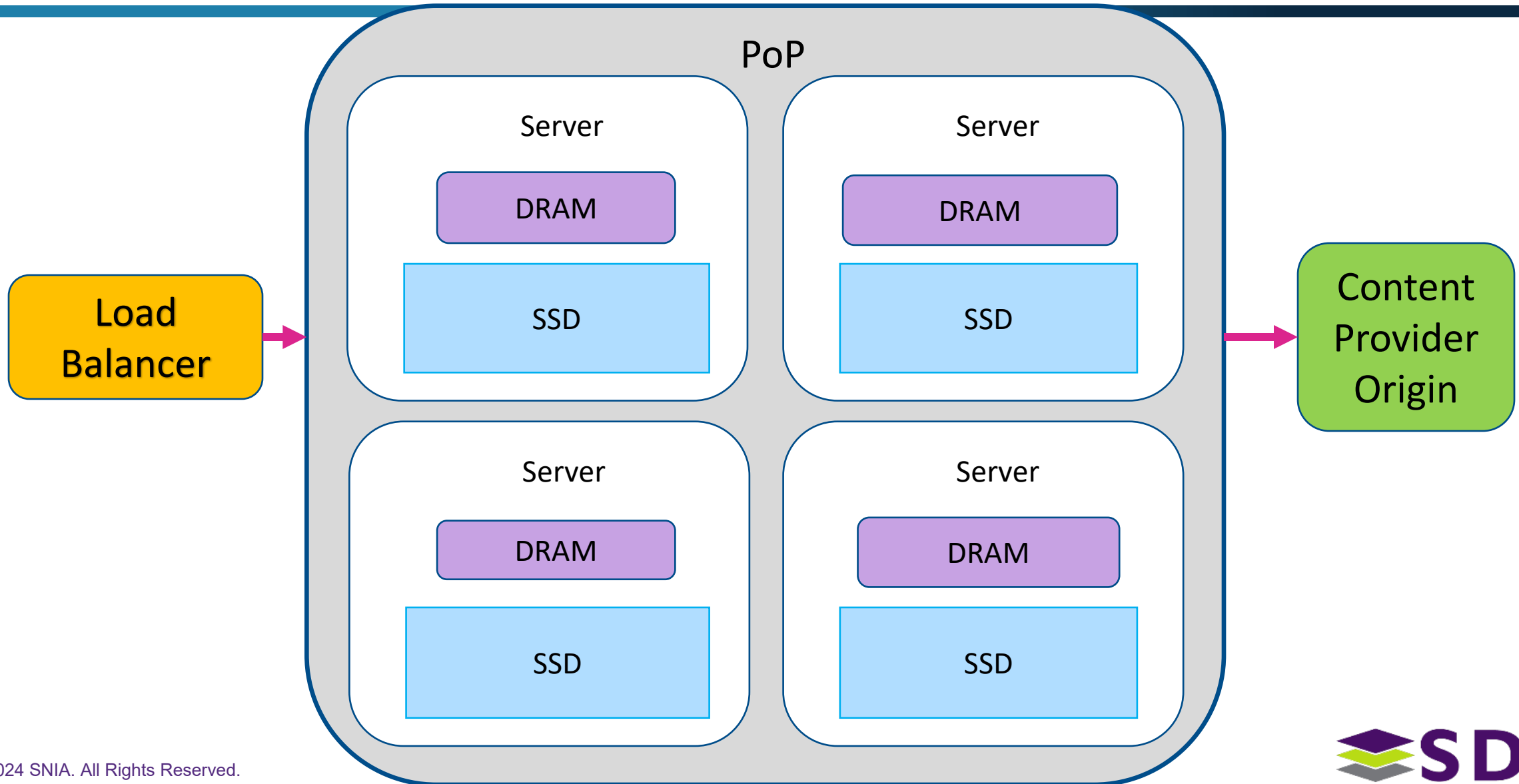
CDN PoP – Compute + Memory + Storage



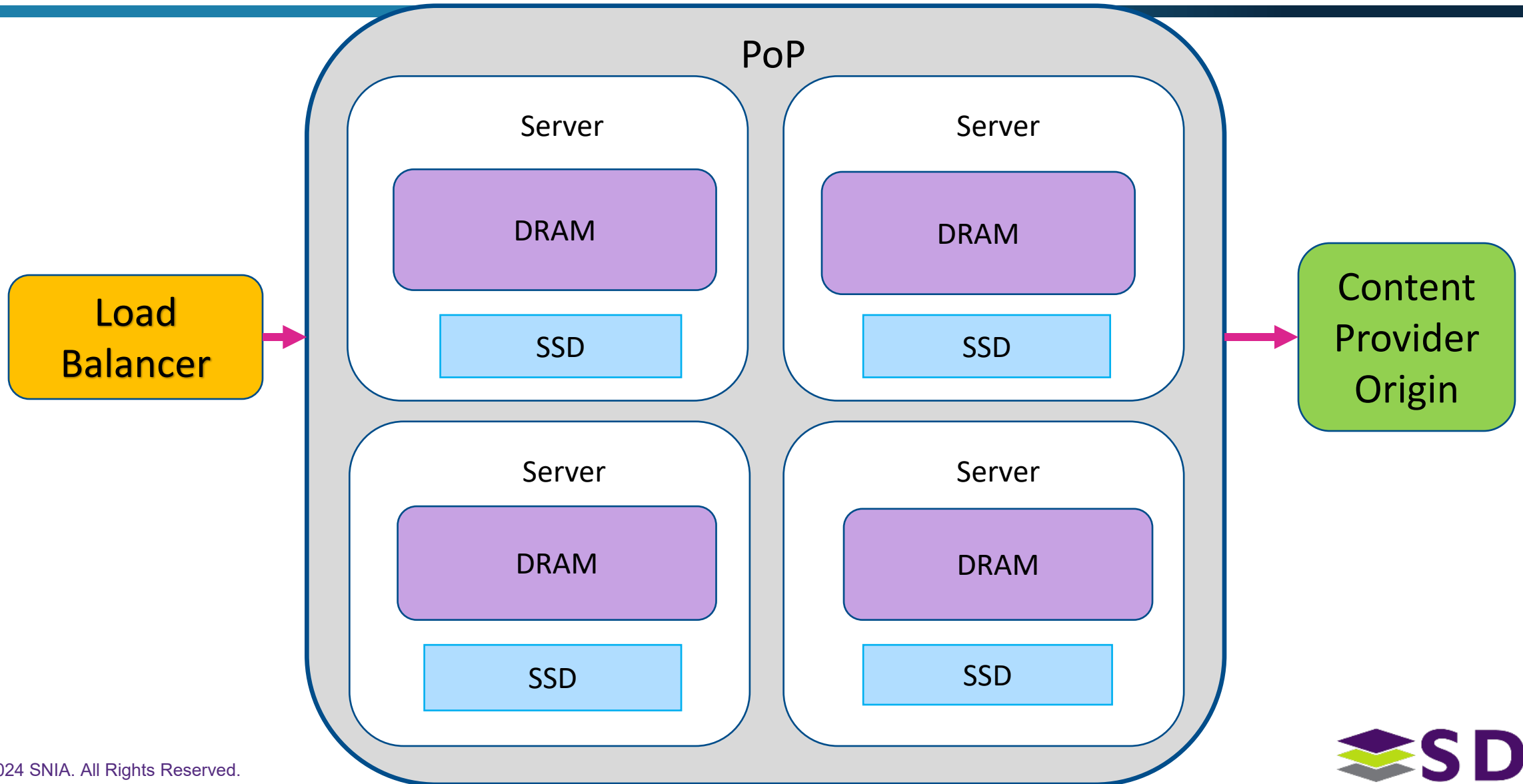
CDN PoP – Compute + Memory + Storage



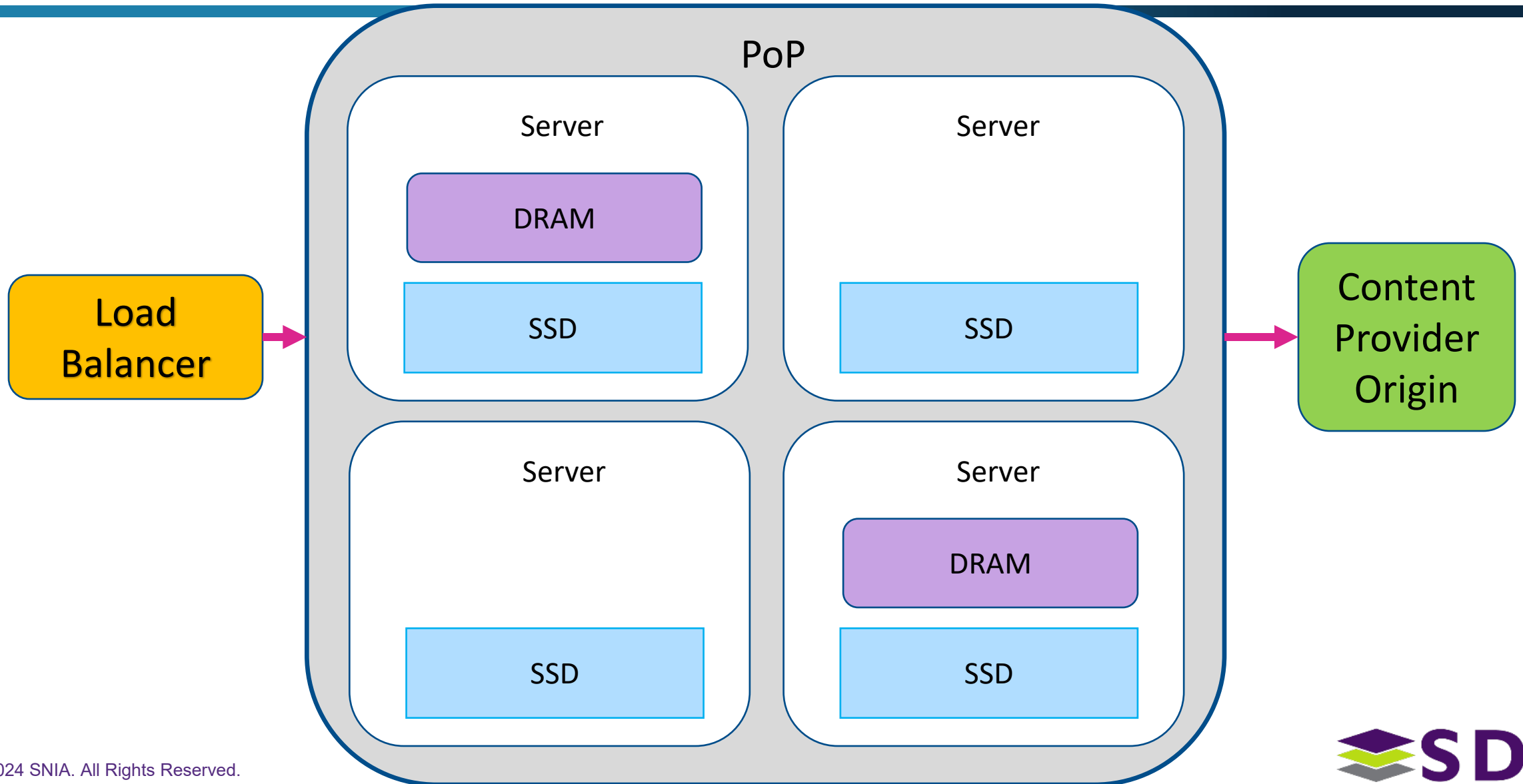
CDN PoP – Compute + Memory + Storage



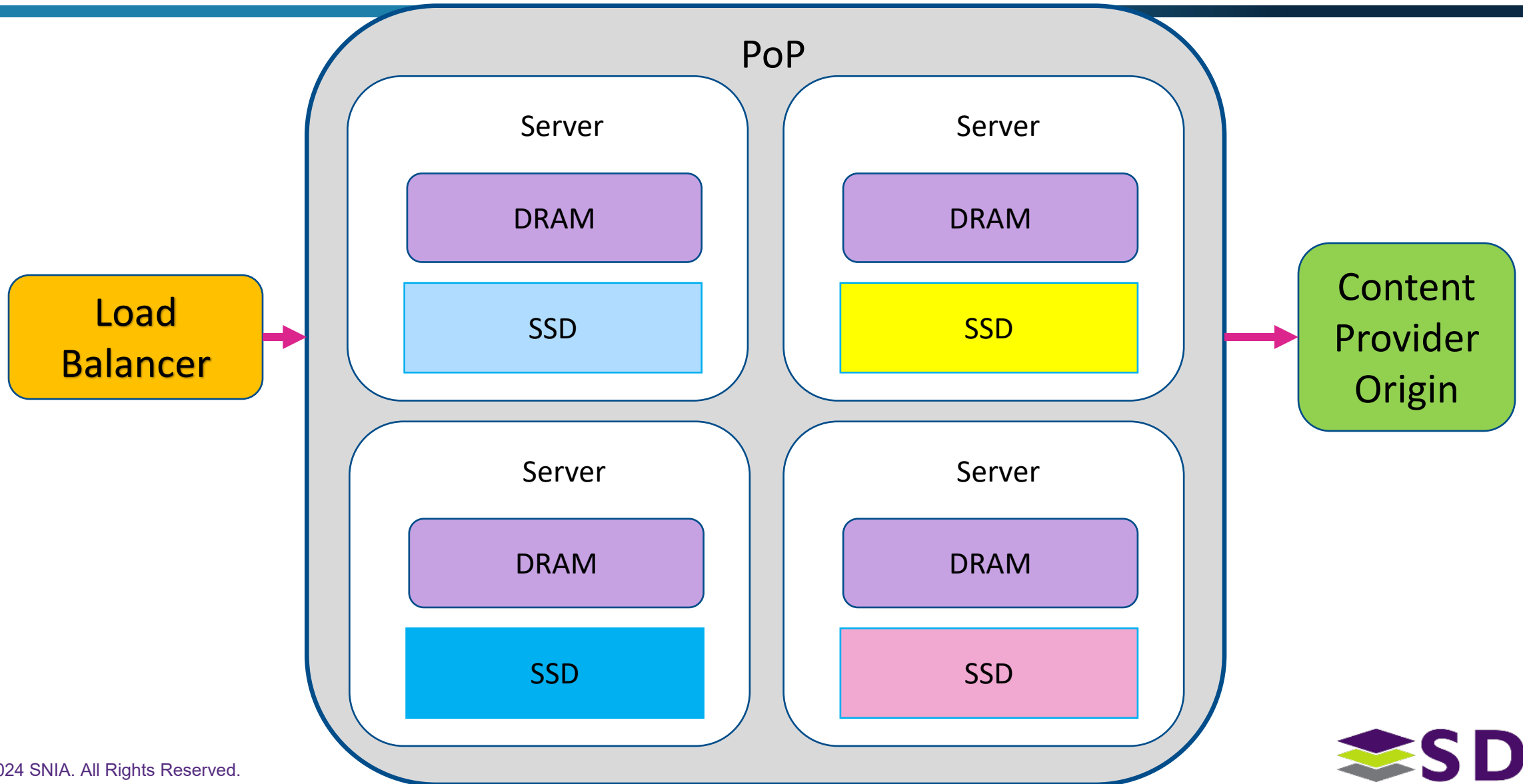
CDN PoP – Compute + Memory + Storage



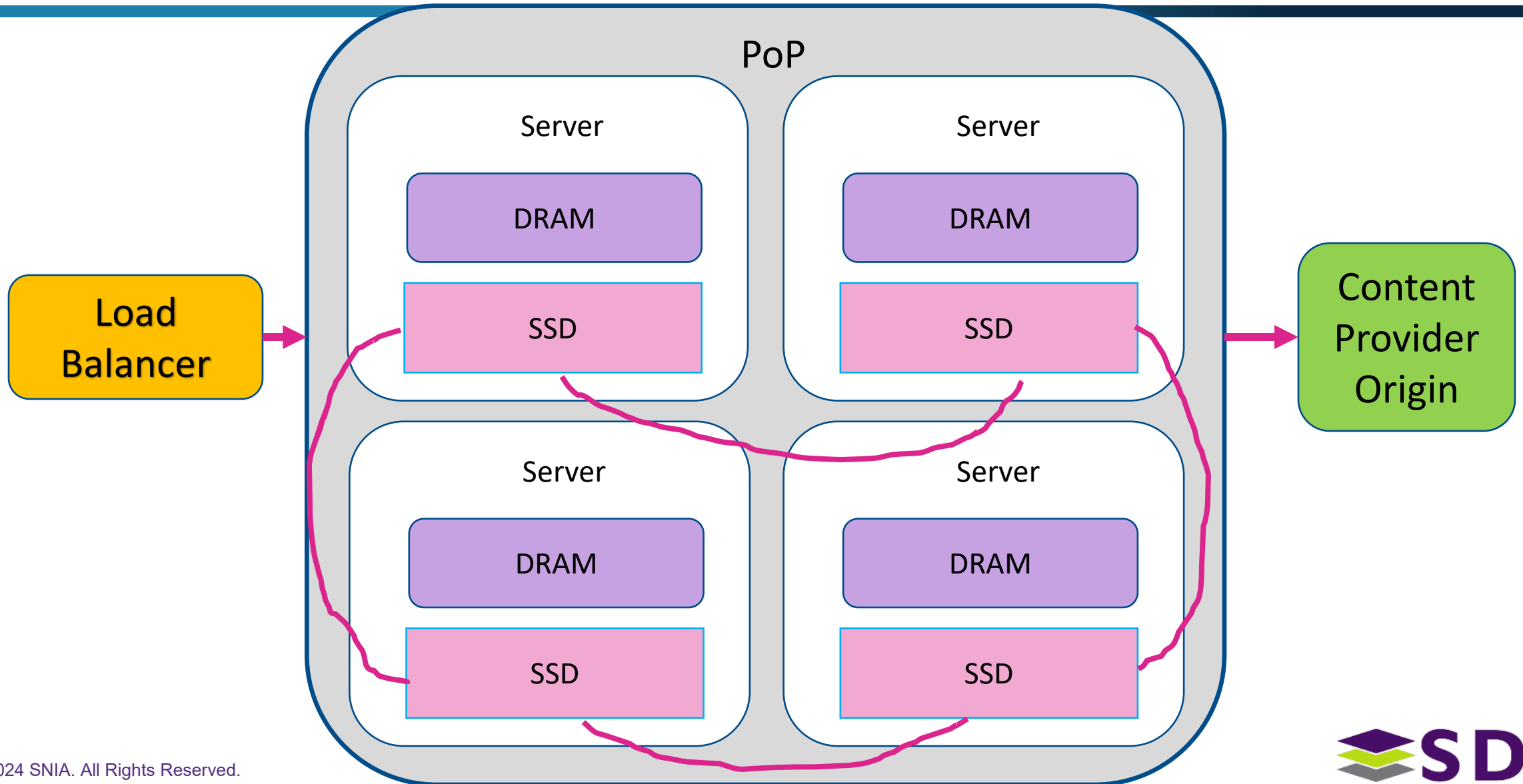
CDN PoP – Compute + Memory + Storage



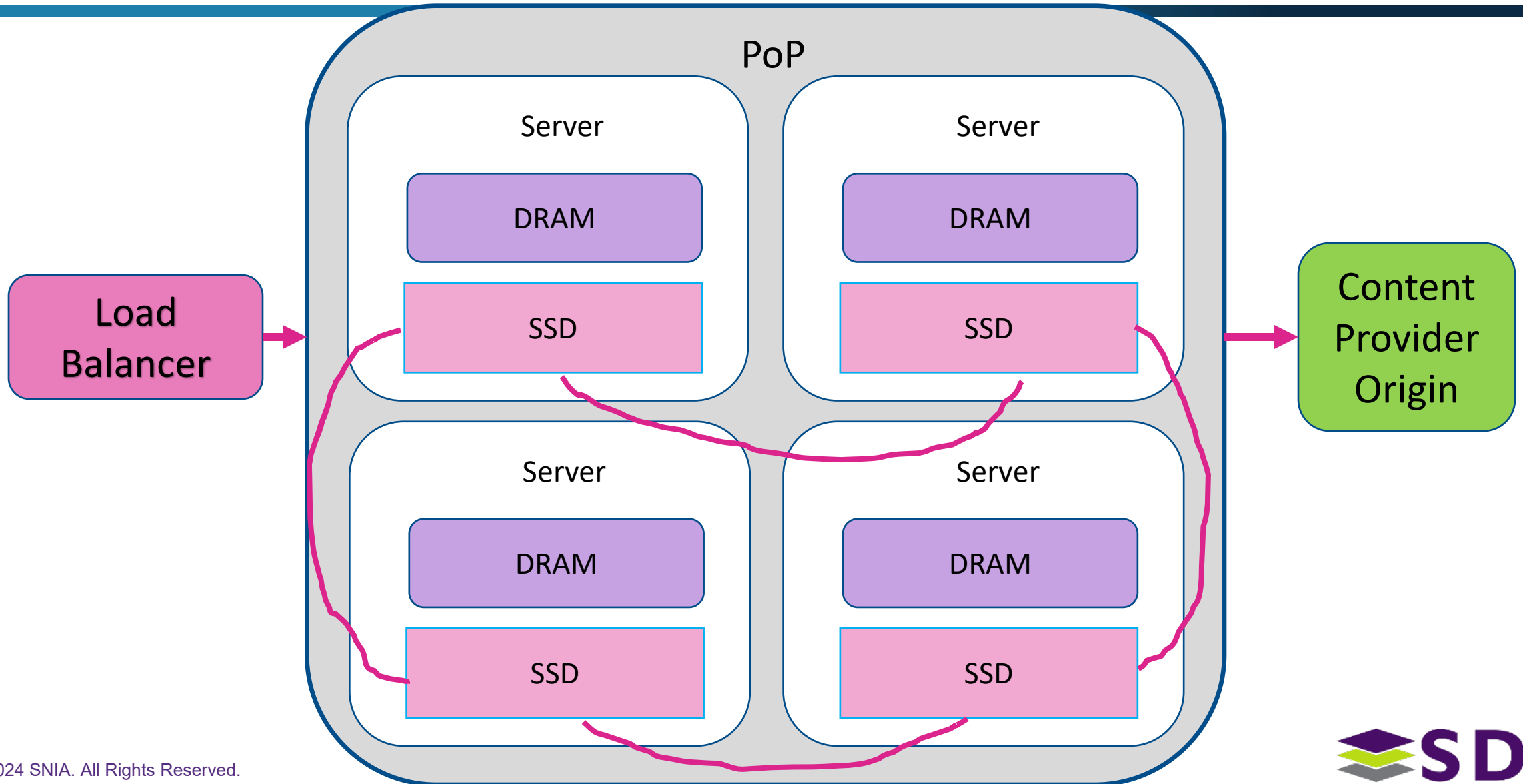
CDN PoP – Compute + Memory + Storage



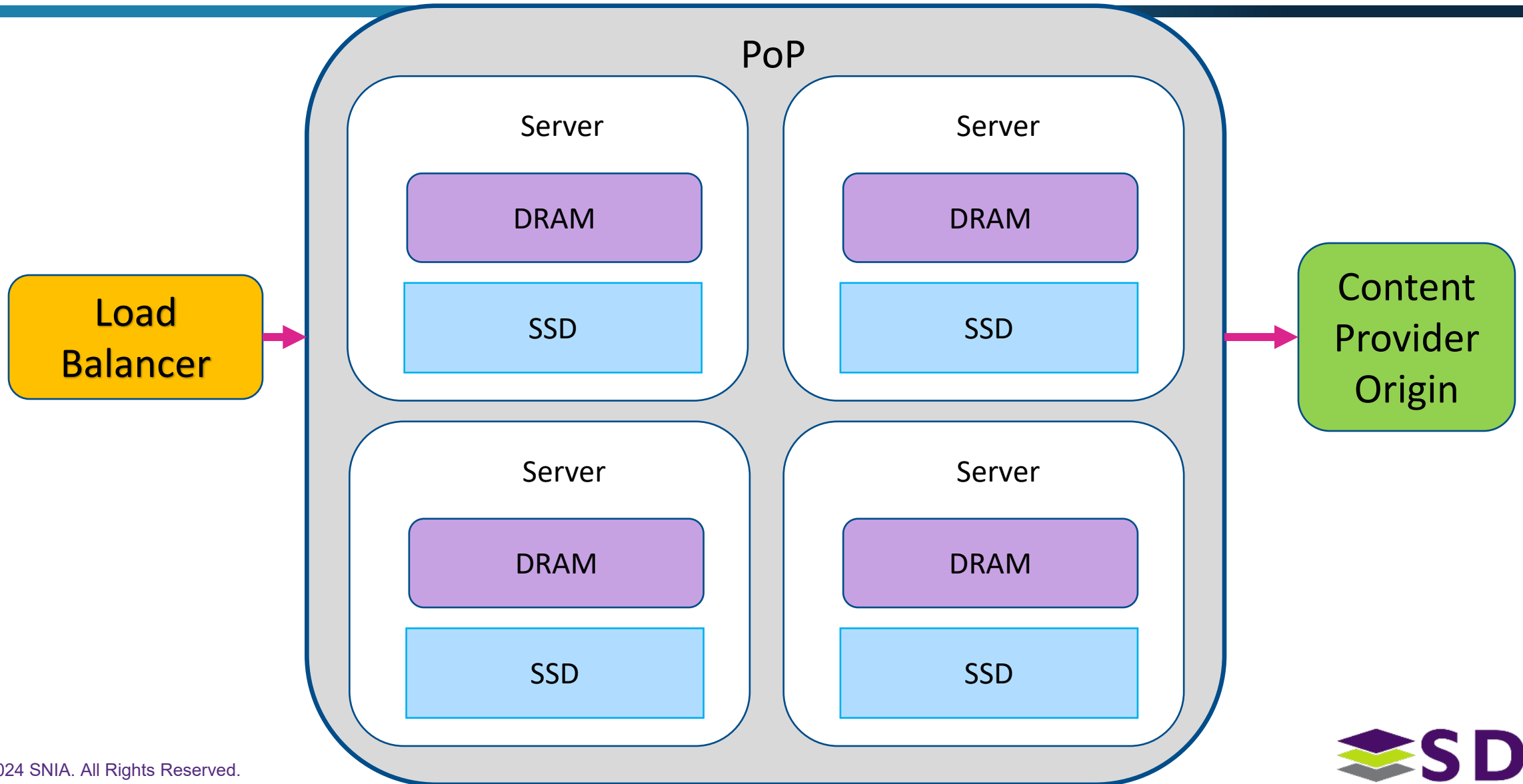
CDN PoP – Compute + Memory + Storage



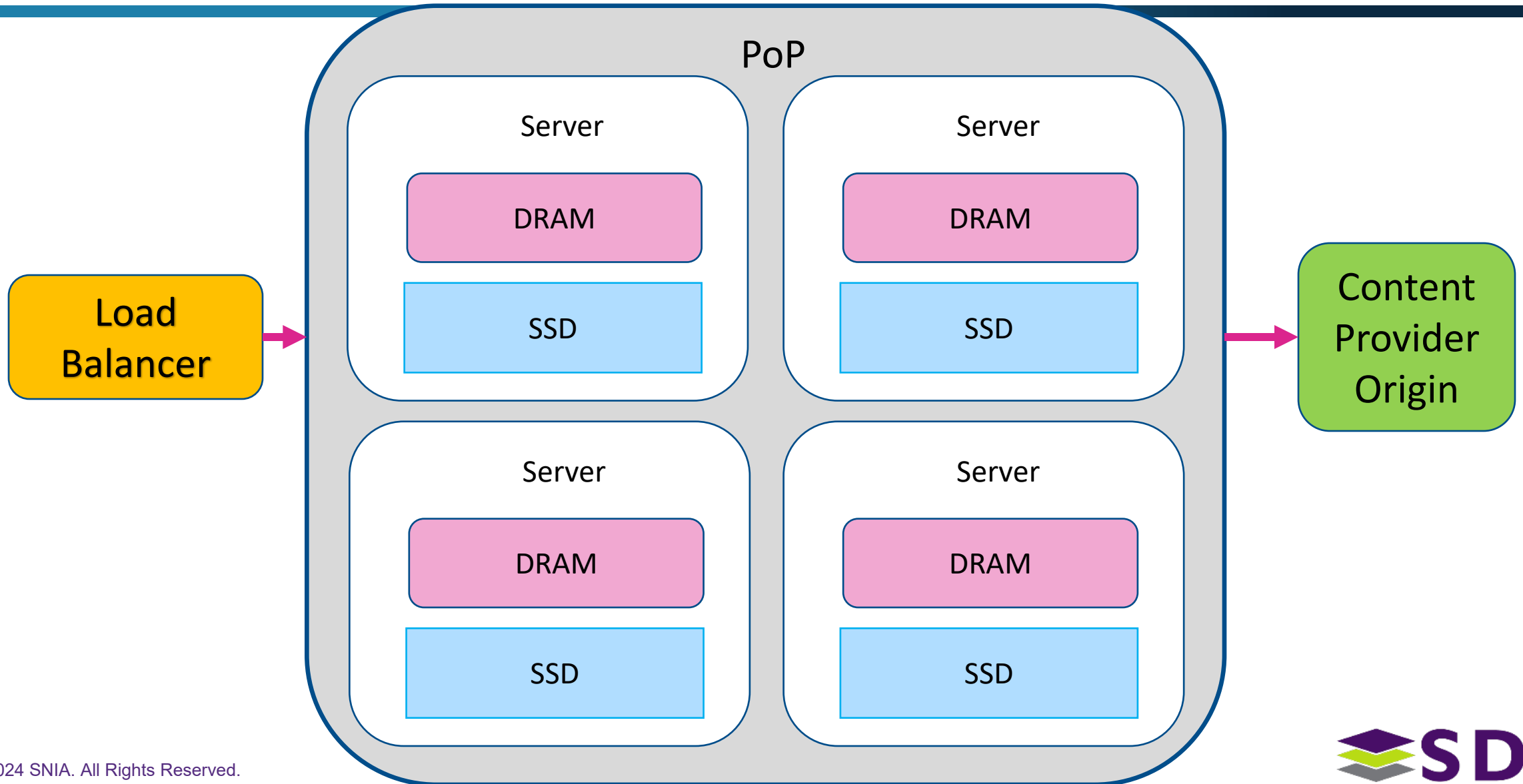
CDN PoP – Compute + Memory + Storage



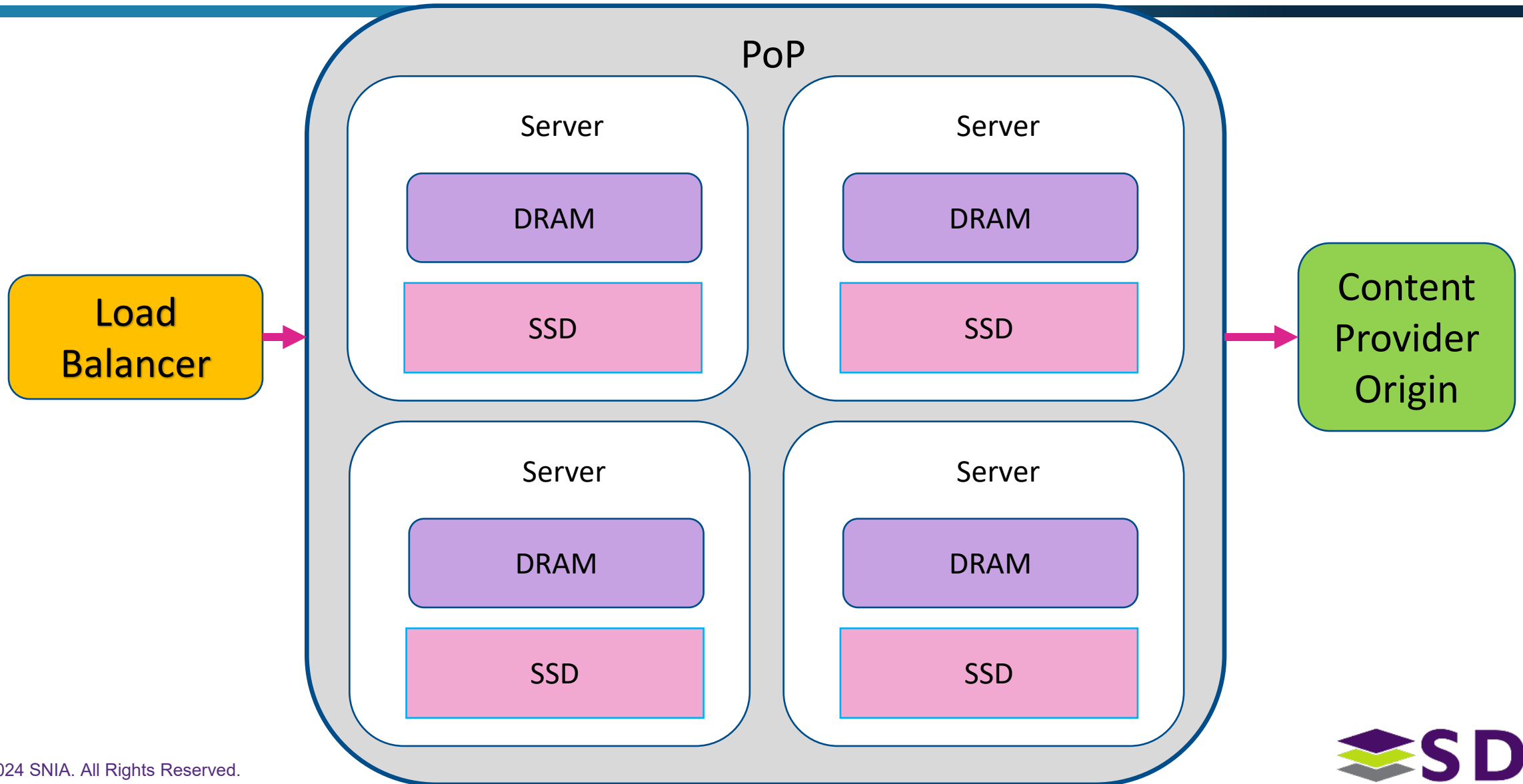
CDN PoP – Compute + Memory + Storage



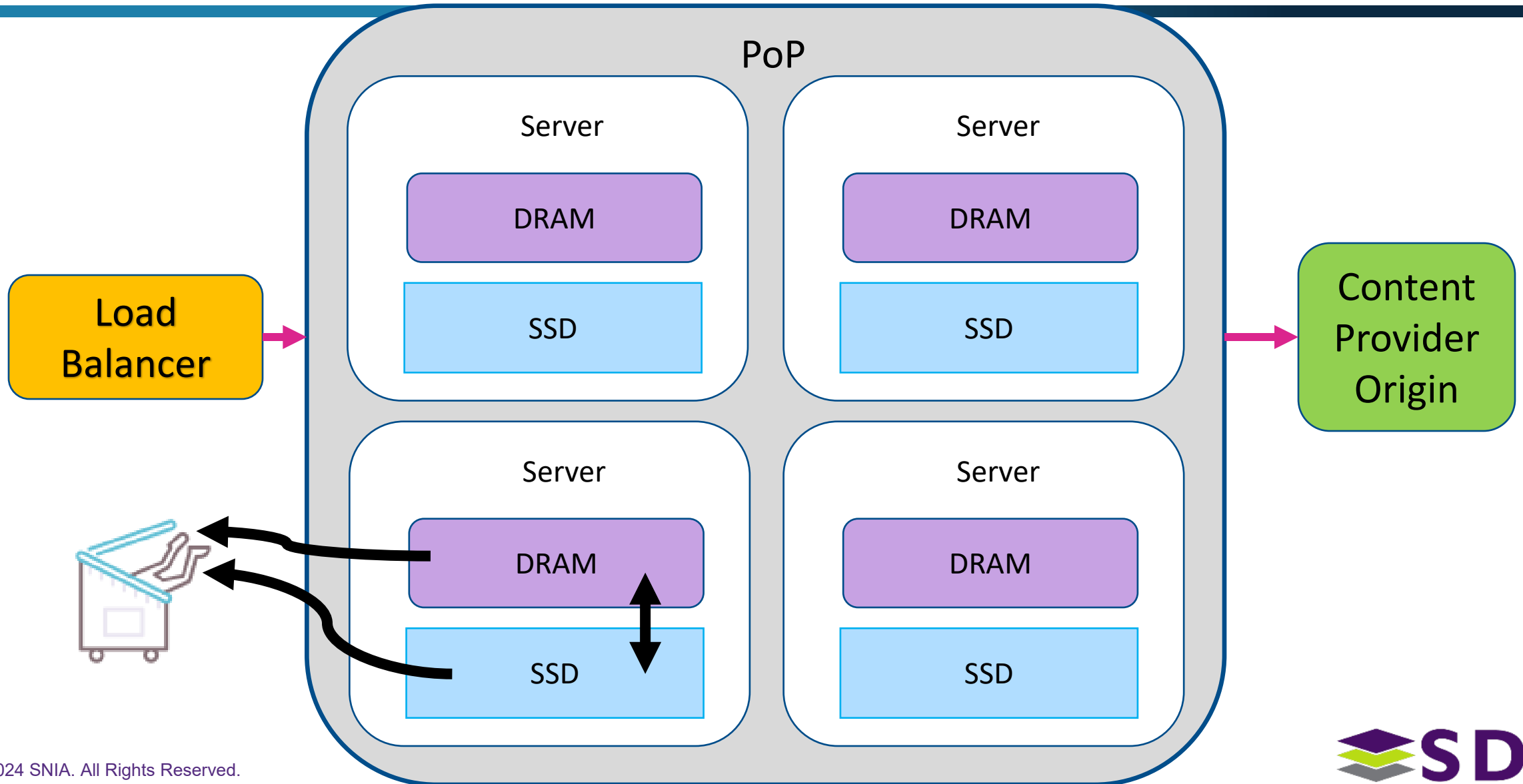
CDN PoP – Compute + Memory + Storage



CDN PoP – Compute + Memory + Storage



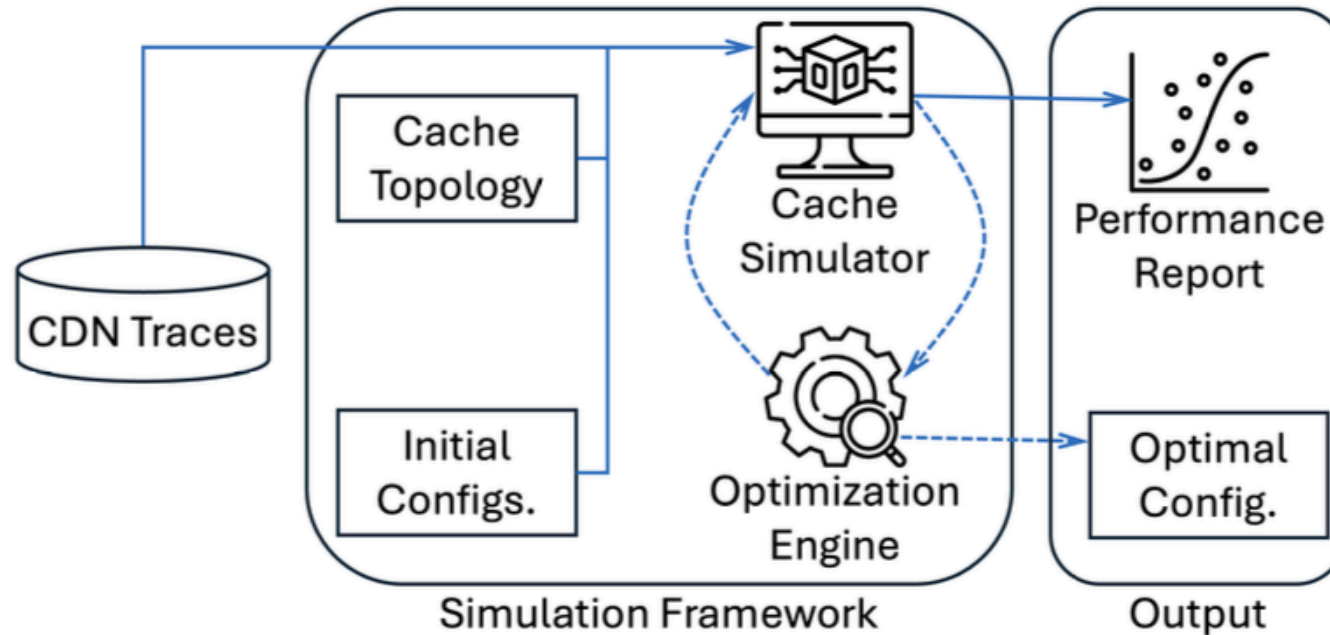
CDN PoP – Compute + Memory + Storage



Case Studies

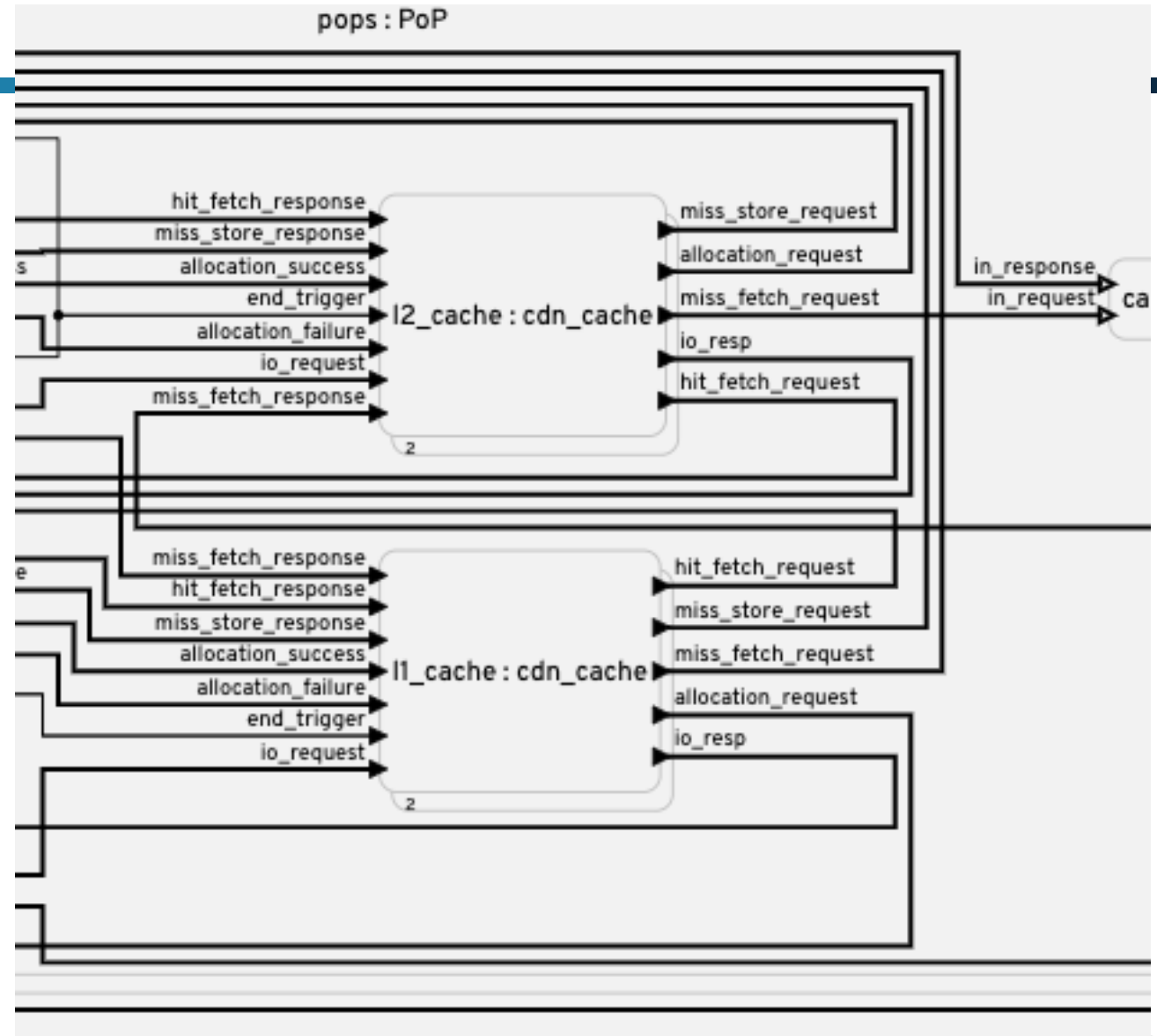
How to Build a Model

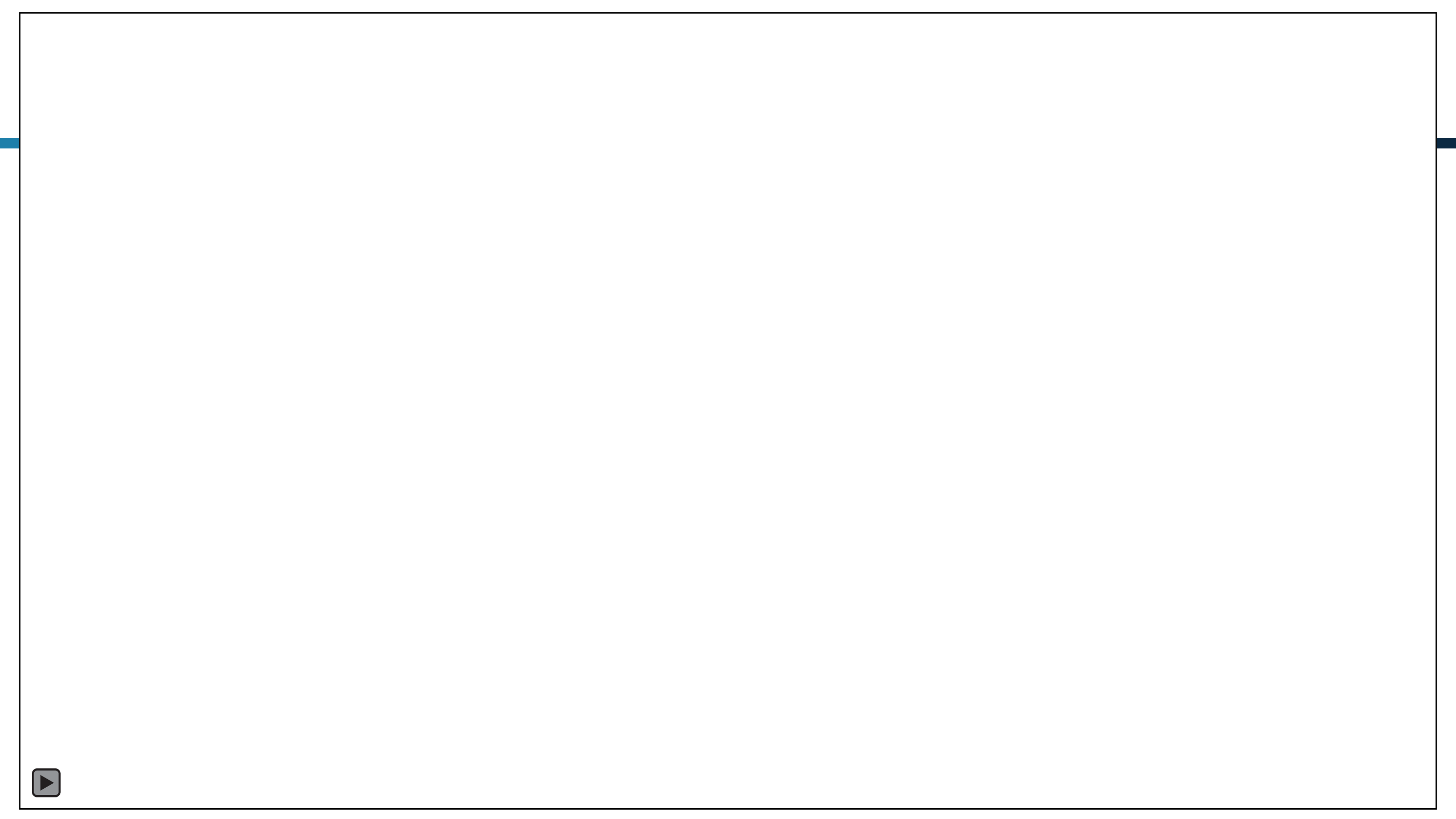
- Conceptual view of the design



Details, Details

- Cheaper, faster and more flexible than hardware models
- Each component can be modeled
- Variables are easy to introduce





Simulation Code

- UI generated from code
- Code simulates component

```
reactor cdn_cache (bank_index:int(0), pop_tier_id:int(0), pop_id:int(0), n_ports:int(1), cache_level:string("L1"), cache_size:uint32_t(4096), page_size:uint32_t(4096))
{
    input io_request:cdn_cache_request_t;
    output io_resp:cdn_cache_response_t;
    output hit_fetch_request:cdn_cache_request_t;
    input hit_fetch_response:cdn_cache_response_t;
    output miss_fetch_request:cdn_cache_request_t;
    input miss_fetch_response:cdn_cache_response_t;
    output miss_store_request:p_cdn_cache_entry_t;
    input miss_store_response:p_cdn_cache_entry_t;
    output allocation_request:uint32_t;
    input allocation_success:uint64_t;
    input allocation_failure:int;
    input end_trigger:uint32_t;
    logical action sch_response(0):cdn_cache_response_t;
    logical action sch_eviction(0):cdn_cache_request_t;
    logical action sch_insertions(0):cdn_cache_entry_t*;
    state free_space:uint32_t(0);

    CacheCtrl = new controller<int> (cache_level = cache_level, name = "cache_controller", pop_tier_id = pop_tier_id, pop_id = pop_id, cache_id = bank_index, log_level = LOG_DEBUG_LEVEL);
    LookUp = new lookup<cdn_cache_request_t, p_cdn_cache_entry_t, cdn_response_tuple_t> (n_ports = 1, log_level = {=LOG_DEBUG_LEVEL=});
    Eviction = new eviction<p_cdn_cache_entry_t, cdn_cache_request_t, cdn_response_tuple_t> (evict_methods = {=&lru_eviction_methods=}, n_ports = 1, eviction_type = LRU);
    Allocator = new allocator<cdn_cache_request_t, p_cdn_cache_entry_t, cdn_response_tuple_t> (log_level = {=LOG_DEBUG_LEVEL=});

    reaction (startup) -> allocation_request {=
        self->free_space = self->cache_size - (self->cache_size % self->page_size);

        LOG_INFO (self->log_level, "(%lld, %u) physical_time:%lld "
            "cdn_cache_%d_%d %s_%d startup paged_cache_size:%u",
            lf_time_logical_elapsed(), lf_tag().microstep, lf_time_physical_elapsed(),
            self->pop_tier_id, self->pop_id, self->cache_level, self->bank_index, self->free_space
        );

        lf_set (allocation_request, self->free_space);
    }
}
```

Case Studies

- Cloudflare and Wikimedia
- Genuine workloads tested
- Variants of baseline algorithms
 - Number of L1 and L2 caches
 - Promotion and demotion policies
 - Eviction policies of L1 and L2
 - DRAM/SSD ratios
 - Load balancer algorithm
- ~125k variants, each run against a day's worth of traffic

Variables

```
79 storage_capacity:
80   value: [549755813888]
81 cache_size_ratio:
82   value: [0.001, 0.005, 0.01, 0.1, 0.2, 0.3, 0.4, 0.5]
83 store_on_miss:
84   value: ["false"]
85 store_from_origin:
86   value: ["true"]
87 move_on_hit_source:
88   value: ["false"]
89 move_on_hit_sink:
90   value: ["false"]
91 move_on_hit_count:
92   value: [1]
93 eviction_types:
94   value: [FIFO, SIEVE, CLOCK, LRU]
95 eviction_methods:
96   value: ["&lru_eviction_methods"]
97 page_size:
98   value: [4096]
99 l2_server_lb:
100   name: L2_SERVER_LOADBALANCER
101 selection_types:
102   value: [USR_DEF_SELECTION]
103 selection_methods:
104   value: ["&url_hash_lb_methods"]
105 l2_servers:
106   value: [1, 2, 3, 4, 5, 6, 7, 8, 9, 10, 30]
107   name: L2_SERVER
108   switch:
109     value: [enable]
110   storage_media:
111     value: [ssd_lite]
112   storage_capacity:
113     value: [549755813888]
114 store_on_miss:
115   value: ["false"]
116 store_from_origin:
117   value: ["true"]
118 move_on_hit_source:
119   value: ["false"]
120 move_on_hit_sink:
121   value: ["false"]
122 eviction_types:
123   value: [FIFO, SIEVE, CLOCK, LRU]
124 eviction_methods:
125   value: ["&lru_eviction_methods"]
126 page_size:
127   value: [4096]
128 pop_serializer:
129   name: POP_SERIALIZER
130 data_mover:
131   name: POP_L1_TO_L2_MOVER
132   switch:
133     value: [disable]
```

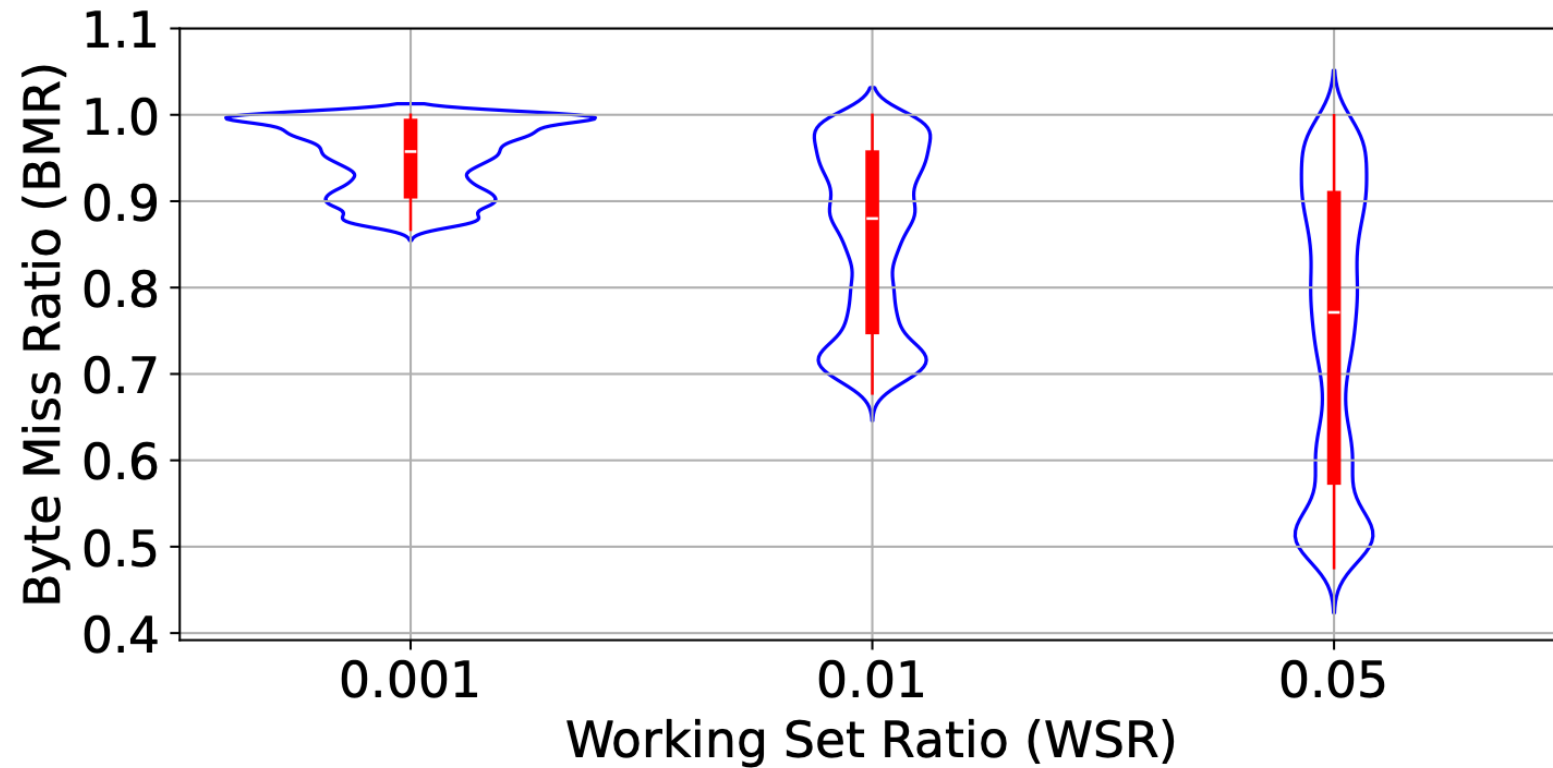
Wiki
Topology

```
79 storage_capacity:
80   value: [549755813888]
81 cache_size_ratio:
82   value: [0.001, 0.005, 0.01, 0.1, 0.2, 0.3, 0.4, 0.5]
83 store_on_miss:
84   value: ["false"]
85 store_from_origin:
86   value: ["true"]
87 move_on_hit_source:
88   value: ["true"]
89 move_on_hit_sink:
90   value: ["false"]
91 move_on_hit_count:
92   value: [1, 2, 3, 4]
93 eviction_types:
94   value: [FIFO, SIEVE, CLOCK, LRU]
95 eviction_methods:
96   value: ["&lru_eviction_methods"]
97 page_size:
98   value: [4096]
99 l2_server_lb:
100   name: L2_SERVER_LOADBALANCER
101 selection_types:
102   value: [LB_1_ON_1_WIRING_DEFAULT]
103 selection_methods:
104   value: ["&url_hash_lb_methods"]
105 l2_servers:
106   value: [1, 2, 3, 4, 5, 6, 7, 8, 9, 10, 30]
107   name: L2_SERVER
108   switch:
109     value: [enable]
110   storage_media:
111     value: [ssd_lite]
112   storage_capacity:
113     value: [549755813888]
114 store_on_miss:
115   value: ["false"]
116 store_from_origin:
117   value: ["false"]
118 move_on_hit_source:
119   value: ["false"]
120 move_on_hit_sink:
121   value: ["true"]
122 eviction_types:
123   value: [FIFO, SIEVE, CLOCK, LRU]
124 eviction_methods:
125   value: ["&lru_eviction_methods"]
126 page_size:
127   value: [4096]
128 pop_serializer:
129   name: POP_SERIALIZER
130 data_mover:
131   name: POP_L1_TO_L2_MOVER
132   switch:
133     value: [enable]
```

Cloudflare
Topology

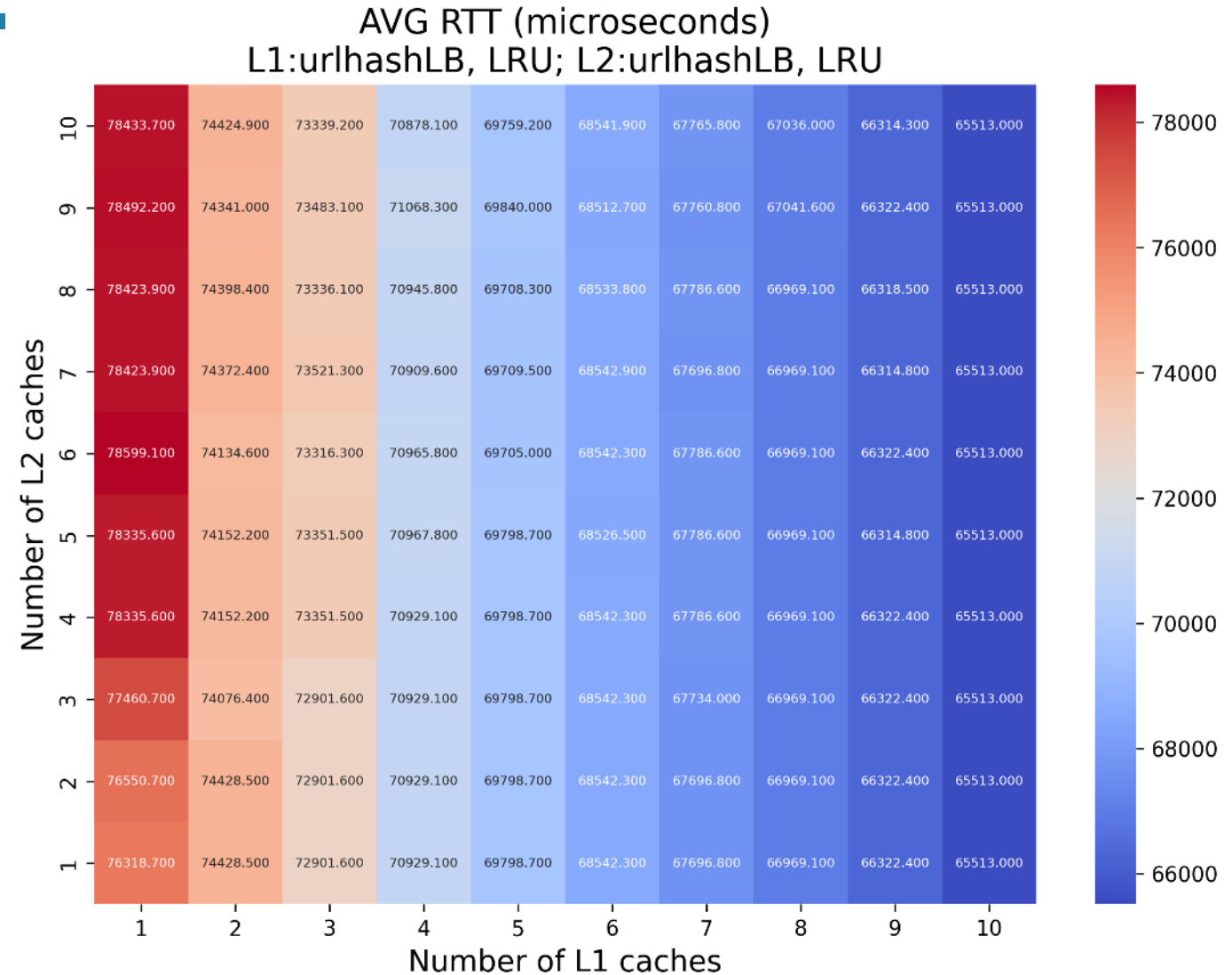
Lots of Variability, Lots of Results

- Multi-objective Optimization

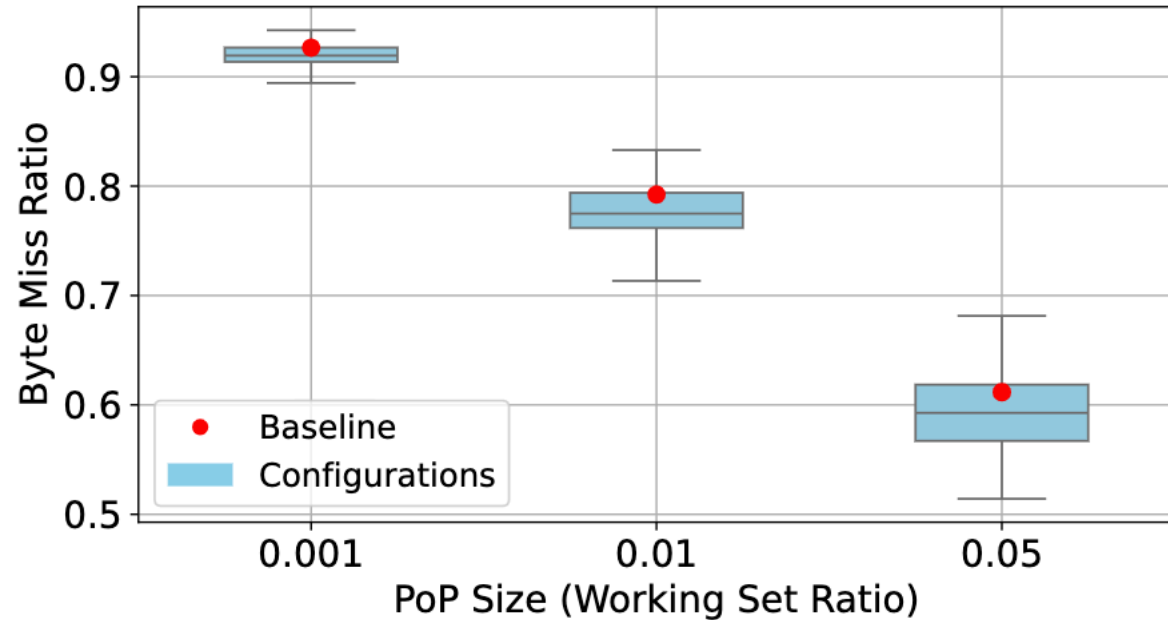


Results Are Easy to Compare

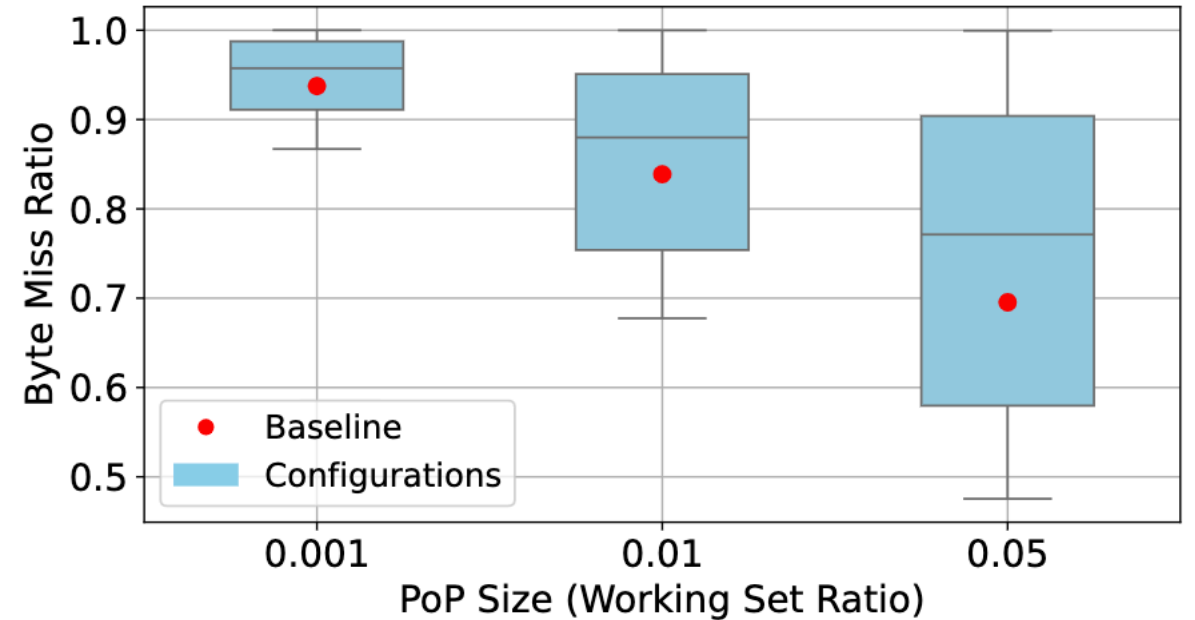
- See many thousands of runs
- More variables = more options



Case Study Results

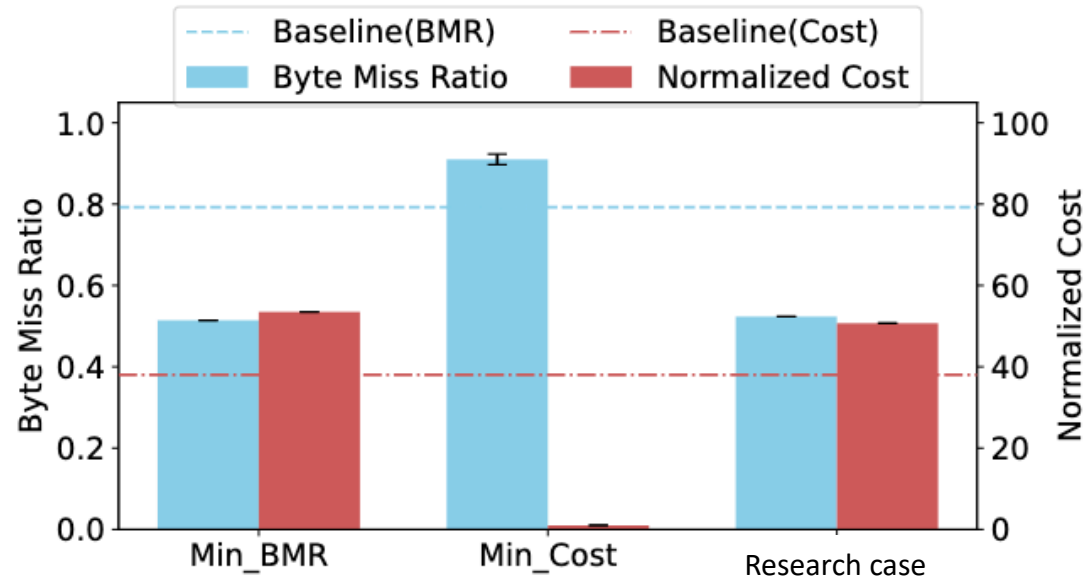


Wikimedia

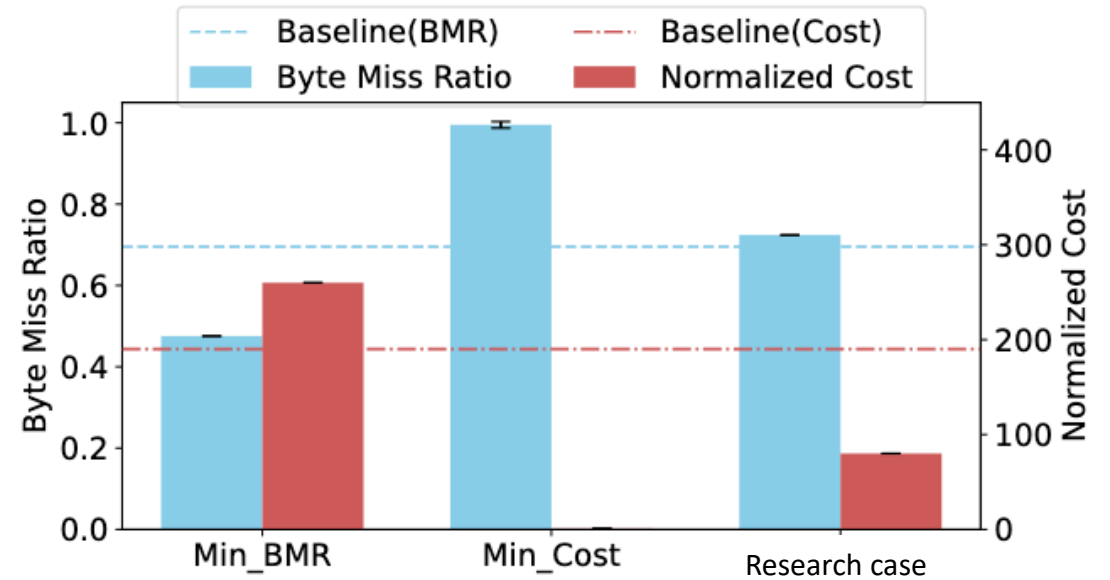


Cloudflare

Case Study Results

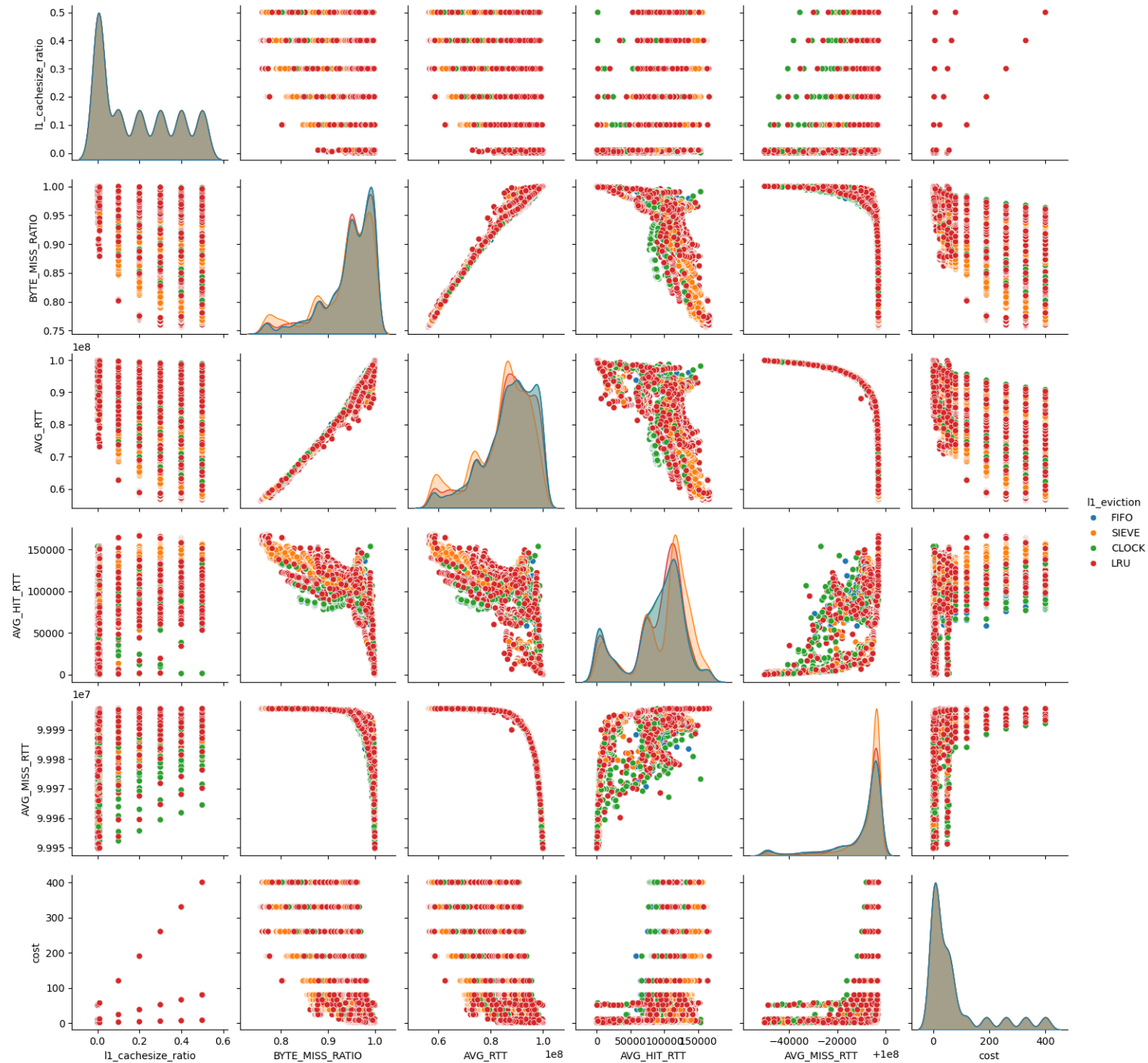


Wikimedia

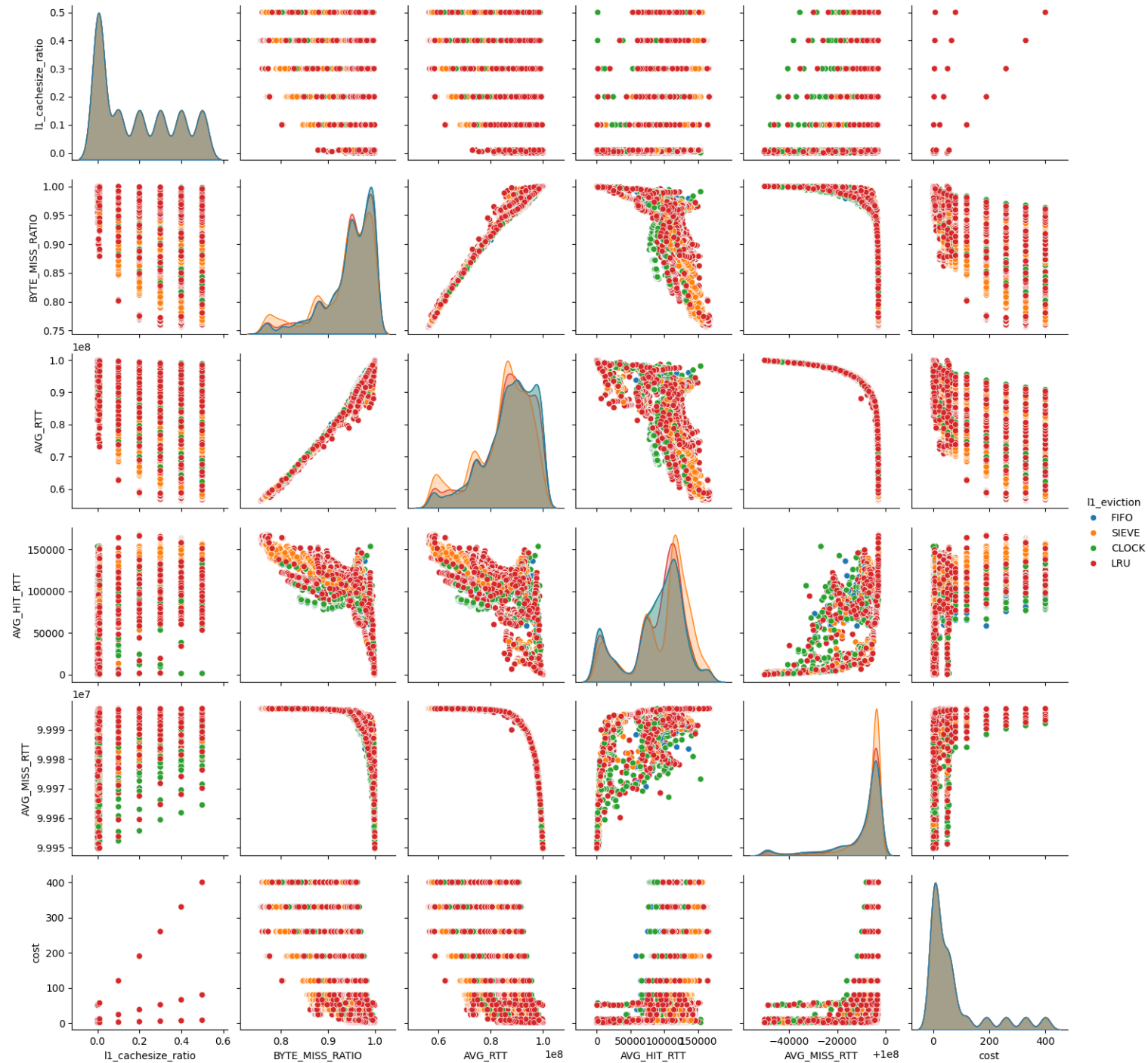


Cloudflare

Results

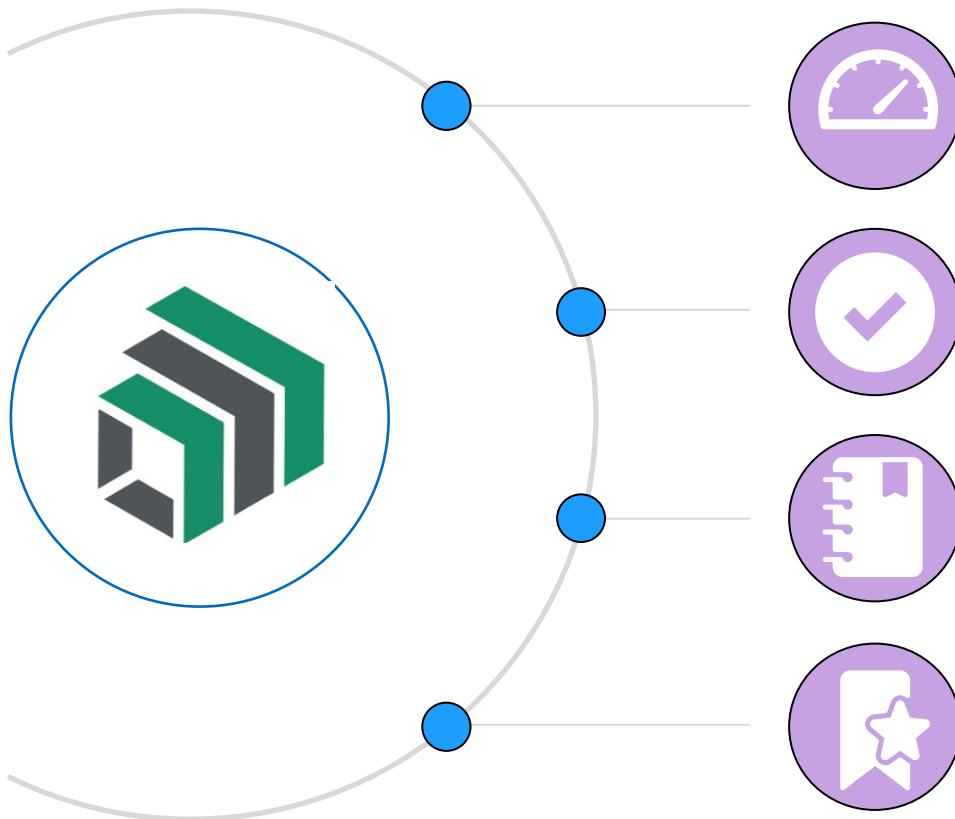


Results



ABOUT MAGNITION

STORAGE PERFORMANCE, REINVENTED



World's First Real-Time Data Placement Optimization
Patented technology is a first for the industry.

Proven At-Scale, with Production Workloads
Use customer traces to fully test diverse workloads in real-time.

Peer-Reviewed and Published in Leading Journals
Multiple industry articles published and reviewed.

Award-Winning, Patented Technology
3-time award winner for innovative technology.

RESULTS WITH MAGNITION

PROVEN IN MARKET TODAY

As an example, a current customer has achieved the following measurable outcomes with Magnition:

Experiments **per day per engineer**

- Without Magnition: **2**
- With Magnition: **50,000+**

Parameter variations tested **before prod release**

- Without Magnition: **50**
- With Magnition: **1,000,000+**

Workload performance improvement using our products to find **optimal out-of-the-box settings: 10-50%+**



I Learned Something Today

1. Caching is used for more than storage
2. Real world CDN implementations can be improved with low effort
3. Large scale simulations can drive huge gains in efficiency and cost



Please take a moment to rate this session.

Your feedback is important to us.