

Cisco Silicon One G200

Contents

Value statement	3
Product overview	3
Features and benefits	4
Prominent feature	5
Product sustainability	9
For more information	9

Value statement

Web-scale, enterprise, and service provider hardware is built around switching silicon, routing line card silicon, and routing fabric silicon. These three basic building blocks enable silicon and system vendors to create unique architectures tuned for individual markets and industries. However, forcing customers to consume and manage these disjointed, dissimilar products has also caused an explosion in complexity, CapEx, and OpEx.

The Cisco Silicon One™ architecture ushers in a new era of networking, enabling one silicon architecture to address a broad market space, while simultaneously providing best-of-breed devices. Cisco Silicon One doesn't mean one device across the network, but one architecture and many optimized devices across the network.

At 51.2 Tbps, the Cisco Silicon One G200 builds on the groundbreaking technology of the Cisco Silicon One G100. What's more, it fully optimizes the design for high-bandwidth, web-scale switching and Artificial Intelligence (AI) and Machine Learning (ML), enabling a deterministic, low-latency, and power-efficient 64x800GE switch.

Product overview

The Cisco Silicon One G200 processor is a 51.2-Tbps, full-duplex, standalone switching processor that can be used to build fixed form factor switches ideally targeted for web-scale data center spine and leaf applications and AI/ML applications.

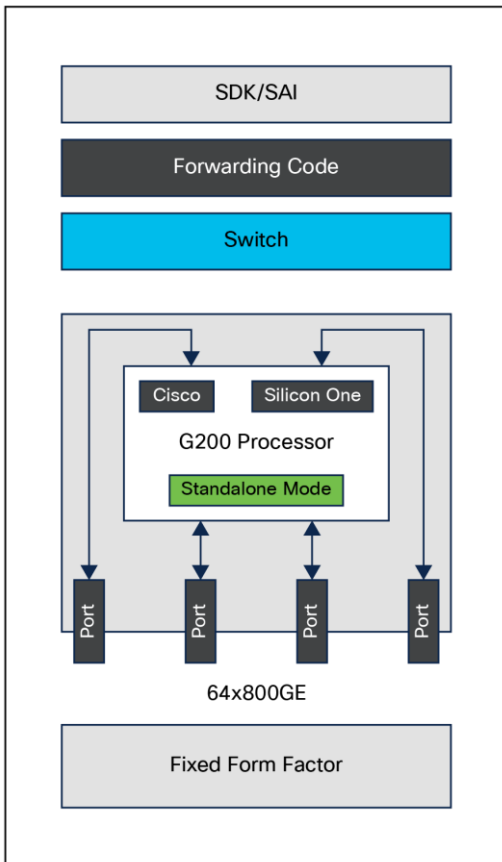


Figure 1.
Form Factor

Features and benefits

Table 1. Architectural characteristics and benefits

Feature	Benefit
One architecture across multiple markets	One architecture, greatly simplifies customer network infrastructure deployments, saving both OpEx and CapEx while simultaneously shortening qualification time.
One SDK across market segments and applications	One SDK provides a consistent point of integration for all applications across the entire network infrastructure, improving quality while reducing OpEx and CapEx for customers.
Latency optimized programmable network processor	Deterministic and low-latency programmable processor that offers additional run-to-completion flexibility for complex flows. This architecture uniquely addresses the requirements of web-scale providers' switching applications without sacrificing features and programmability.
Large and fully unified packet buffer	Fully shared on-die packet buffer allows any input or output port to consume the entire memory. This capability reduces packet loss and minimizes PFC events, thus maximizing network performance under varying traffic conditions and enabling low latency for RDMA and RoCEv2 protocols.
Unmatched telemetry and visibility	Support for standard and emerging web-scale, in-band telemetry protocols enables advanced congestion control and advanced flow tracking with temporal dynamics. Together with in-network trigger events, these capabilities enable post-event analysis in hardware time scales.
Advanced load balancing capabilities	Support for stateless and stateful congestion-aware load balancing techniques ensures optimal delivery of packets through the network. This helps to ensure optimal Flow Completion Time (FCT) for traditional web-scale networks and optimal Job Completion Time (JCT) for massive-scale AI/ML networks.
Network resiliency assurance	Support for hardware-based link monitoring and rebalancing of traffic helps ensure optimized network utilization even under link failure conditions in large-scale networks.

Prominent feature

Flexibility and performance for next-generation, web-scale, front-end, and AI/ML networks

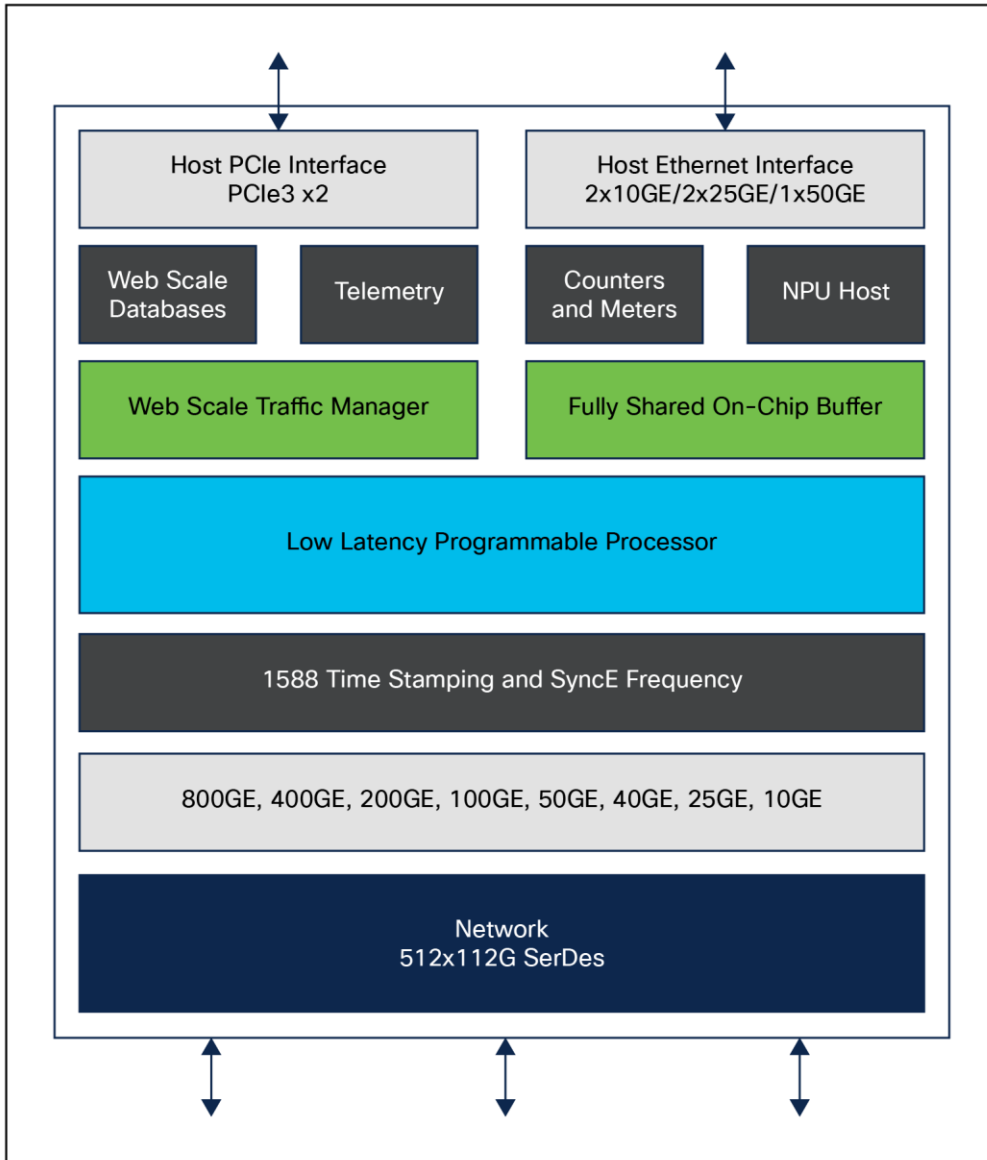


Figure 2.
Block Diagram

Features

- 512 112G long-reach SerDes supporting NRZ and PAM4 modulation
- 512 Ethernet MACs enable maximum network scale-out for optimized network deployments
- Flexible port configuration supporting 10/25/40/50/100/200/400/800 Gbps
- Large, fully shared, on-die packet buffer
- 1588v2 and SyncE support with nanosecond-level accuracy
- On-chip, high-performance, programmable host Network Processing Unit (NPU) for high-bandwidth offline packet processing (for example, OAM processing, MAC learning)
- Multiple embedded processors for CPU offloading
- PCIe gen3 and two Ethernet interfaces to connect to the host CPU complex
- Advanced features for AI/ML deployments with optimal load balancing, fault detection, recovery, and telemetry

Traffic management

- Multiple output queues per output port support web-scale customers' future needs
- Support for ingress and egress traffic mirroring
- Support for link-level (IEEE802.3x) flow control
- Support for PFC priority-level (802.1Qbb) flow control
- Support for PFC watchdog
- Dynamic thresholds and policies help ensure optimal usage of the fully shared packet buffer
- Support for probabilistic multicolor ECN marking
- Support for probabilistic multicolor WRED drop profiles

NPU

- Optimized, deterministic, and low-latency programmable processor with extensions to complete run-to-completion programmable processors for advanced features
- Web-scale optimized and fungible tables
- Achieves line rate at small packet sizes and full web-scale feature sets running

Load balancing

- Flow load balancing using WECMP, ECMP, or LAG with innovative noncorrelated hashing functions to avoid polarization, even across massive-scale networks
- Support for WECMP without replicating entries
- Congestion-aware flowlet load balancing with ability to detect and handle elephant flows
- Congestion-aware packet spraying independent of flow characteristics

Instrumentation and telemetry

- Support for standard (P4-INT, IFA1.0, IFA2.0) and emerging web-scale, in-band telemetry protocols
- Advanced flow scope with temporal dynamics and live network trigger for post-mortem analysis
- Programmable meters used for traffic policing and coloring
- Programmable counters used for flow statistics and OAM loss measurements
- Counters for port utilization, microburst detection, delay measurements, flow tracking, elephant flow detection, and congestion tracking
- Traffic mirroring: (ER)SPAN on congestion and drop
- Support for sFlow and NetFlow

AI/ML

- Advanced load balancing techniques for optimal network performance
- Hardware-based link failure isolation and rerouting to enable performance across large-scale networks
- Advanced congestion control and telemetry to optimize Job Completion Time (JCT)
- Deterministic low-latency performance

Software

- SDK APIs in C++, and in C
- Switch Abstraction Interface (SAI)
- Functional simulation environment
- SONiC reference on functional simulator and hardware platform
- Support for x86 and ARM host CPU complexes
- Distribution-independent Linux packaging
- Debug support: gRPC-based CLI and Python shell

Programmability

- Application development is handled by an IDE programming environment
- At compilation, the forwarding application generates low-level register/memory access APIs and higher-level SDK application APIs
- Provides application support for a wide range of data center, service provider, and enterprise protocols
- Ability to develop the SDK and applications running over the SDK over a simulated Cisco Silicon One device

Cisco application

Utilizing Silicon One's extensible programming toolkit, we are always adding features to address new markets and new customer requirements. A sample of features supported includes:

- | | |
|---|---|
| <ul style="list-style-type: none">• IPv4/v6• MPLS• Ethernet Switching<ul style="list-style-type: none">◦ 802.1d, 802.1p, 802.1q, 802.1ad• IP Tunneling<ul style="list-style-type: none">◦ IPinIP◦ GRE◦ VXLAN• Segment Routing<ul style="list-style-type: none">◦ SRv6 uSID◦ MPLS• RDMA Support<ul style="list-style-type: none">◦ Priority Flow Control (PFC) 802.1Qbb◦ Flow Control (802.3x)◦ Probabilistic multicolor ECN marking◦ Probabilistic multicolor WRED drop• Integrated Routing and Bridging (IRB)• HSRP/VRRP• Policy-Based Routing• Security and QoS ACLs• ECMP and LAG (802.3ad)• Multicast<ul style="list-style-type: none">◦ IGMP• Protection (Link/Node/Path and TI-LFA) | <ul style="list-style-type: none">• QoS Classification and Marking• Congestion Management• Telemetry<ul style="list-style-type: none">◦ NetFlow, sFlow◦ In-band Telemetry (P4-INT, IFA, and emerging protocols)◦ (ER)SPAN◦ Packet Mirroring with Appended Metadata◦ Lawful Intercept• Warmboot• DDoS Mitigation<ul style="list-style-type: none">◦ Control-Plane Policing◦ BGP Flowspec• Timing and Frequency Synchronization<ul style="list-style-type: none">◦ SyncE◦ 1588 |
|---|---|

Product sustainability

Information about Cisco's Environmental, Social, and Governance (ESG) initiatives and performance is provided in Cisco's CSR and sustainability [reporting](#).

Table 2. Cisco environmental sustainability information

Sustainability topic		Reference
General	Information on product-material-content laws and regulations	Materials
	Information on electronic waste laws and regulations, including our products, batteries, and packaging	WEEE Compliance
	Information on product takeback and reuse program	Cisco Takeback and Reuse Program
	Sustainability inquiries	Contact: csr_inquiries@cisco.com
Material	Product packaging weight and materials	Contact: environment@cisco.com

For more information

[Learn more](#) about Cisco Silicon One.

Americas Headquarters
Cisco Systems, Inc.
San Jose, CA

Asia Pacific Headquarters
Cisco Systems (USA) Pte. Ltd.
Singapore

Europe Headquarters
Cisco Systems International BV Amsterdam,
The Netherlands

Cisco has more than 200 offices worldwide. Addresses, phone numbers, and fax numbers are listed on the Cisco Website at <https://www.cisco.com/go/offices>.

Cisco and the Cisco logo are trademarks or registered trademarks of Cisco and/or its affiliates in the U.S. and other countries. To view a list of Cisco trademarks, go to this URL: <https://www.cisco.com/go/trademarks>. Third-party trademarks mentioned are the property of their respective owners. The use of the word partner does not imply a partnership relationship between Cisco and any other company. (1110R)