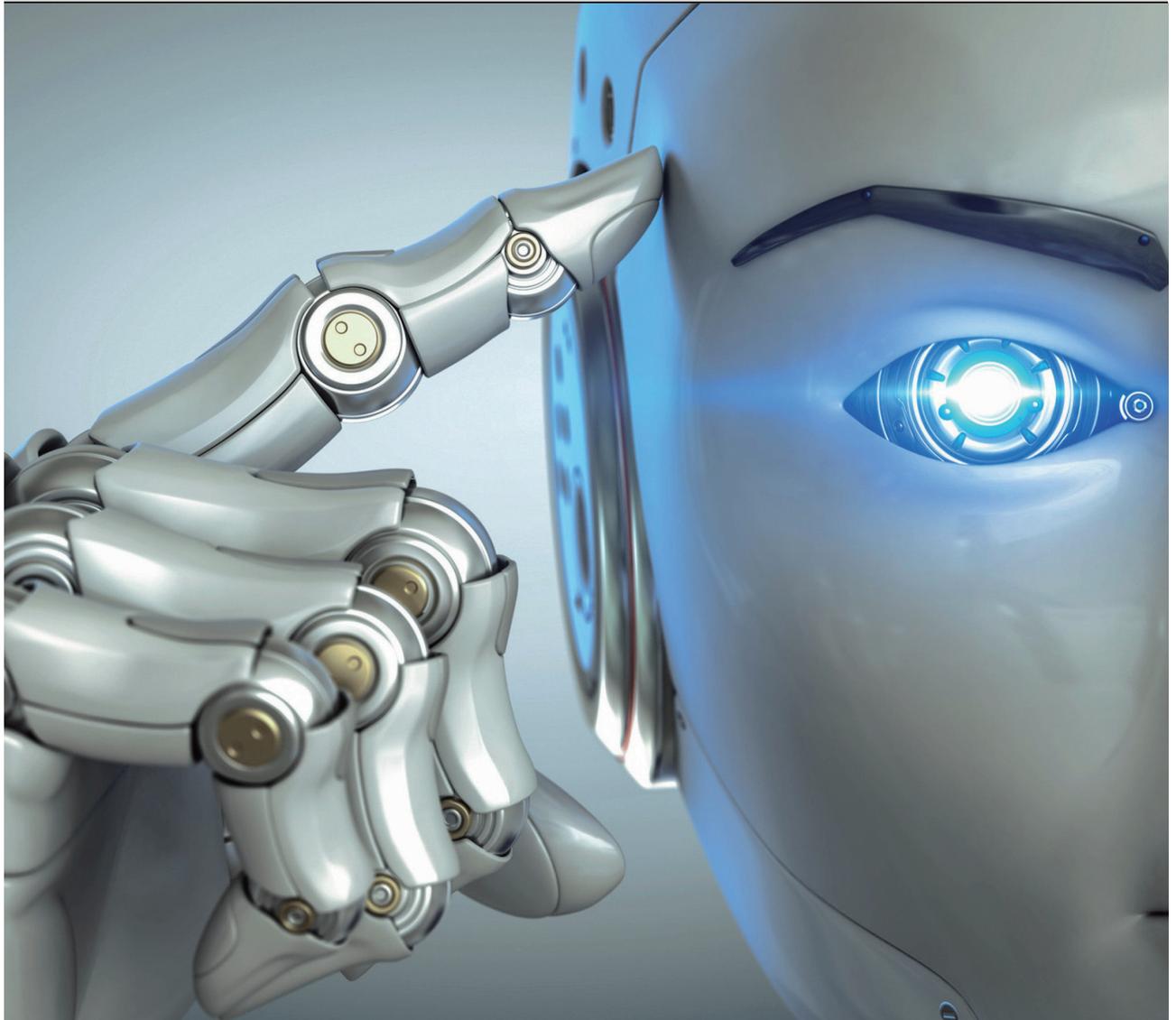


# 전력 소모가 급증하는 컴퓨팅 어플리케이션 솔루션

AI 분야는 오늘날 가장 고도의 컴퓨팅 집약적 과제로 인식되는 작업들을 해결하기 위해 막강한 프로세싱 성능을 요구하고 있다. 이는 온 보드 컴퓨팅 엔진 간의 지연시간을 줄이는 클러스터 아키텍처에 점점 더 높은 성능의 프로세서와 더 큰 메모리 리소스를 구현함으로써 달성되고 있다.

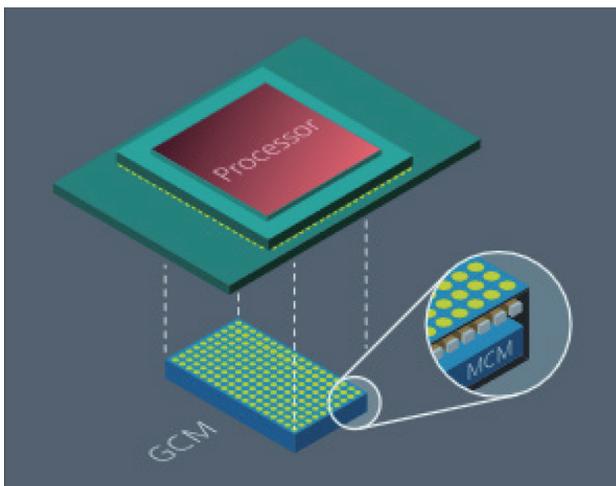
글/로버트 겐드론(Robert Gendron), 바이코의 제품 마케팅 및 기술 리소스 부문 부사장



2019년에는 컴퓨팅 분야의 고객들과 업계의 혁신기  
업들이 소통할 수 있는 수많은 포럼이 진행되었다.  
OCP(Open Compute Project) 글로벌 서밋 및 ODCC  
(Open Data Center Committee)와 같은 데이터센터 중심  
의 이벤트를 비롯해 첫 번째 AI 하드웨어 서밋 및 슈퍼 컴  
퓨팅 2019 이벤트가 개최되었다. 이러한 각 이벤트는 AI와  
슈퍼컴퓨팅 및 클라우드 데이터센터 공급업체들이 프로세  
싱 성능과 전력효율을 극대화하기 위한 핵심 과제에 어떻  
게 접근하고 있는지 확인할 수 있는 자리였다.

AI 분야는 오늘날 가장 고도의 컴퓨팅 집약적 과제로 인  
식되는 작업들을 해결하기 위해 막강한 프로세싱 성능을  
요구하고 있다. 이는 온 보드 컴퓨팅 엔진 간의 지연시간을  
줄이는 클러스터 아키텍처에 점점 더 높은 성능의 프로세  
서와 더 큰 메모리 리소스를 구현함으로써 달성되고 있다.

이러한 상황은 오늘날 AI를 위한 가장 강력한 프로세서  
로 알려진 웨이퍼 스케일 엔진(WSE: Wafer Scale Engine)  
을 최근 발표한 세레브라스(Cerebras)와 같은 회사의 혁신  
을 가속화하고 있다. 웨이퍼 전반에 걸쳐 84개의 프로세싱  
셀로 구성되었지만 단일 칩으로 기능하는 WSE는 기존의  
단일 소켓 기반 칩 아키텍처에서 유발되는 지연시간을 획  
기적으로 줄여준다.



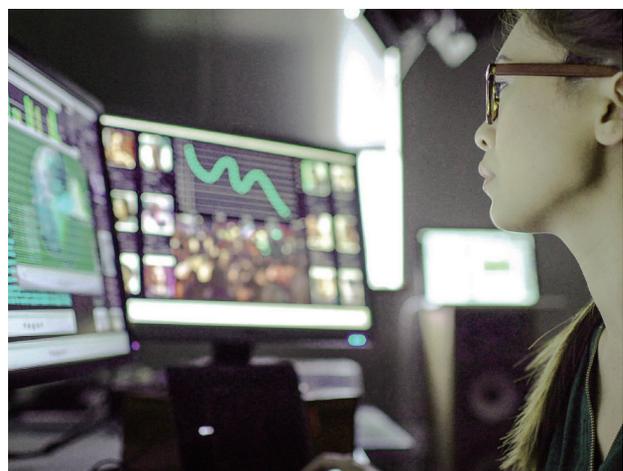
또한 WSE는 기존 프로세서 보다 훨씬 큰 규모의 15kW  
정격을 가지고 있어 매우 높은 전류에서 각 셀로 전력을 균  
일하게 인가할 수 있는 첨단 전력 아키텍처가 요구된다.

이를 달성하기 위해 세레브라스는 VPD(Vertical Power  
Delivery) 아키텍처를 구현한 바이코와 협력하고 있다.  
이 아키텍처는 전력분배 네트워크(PDN: Power Delivery  
Network) 저항을 50% 이상 감소시키는 VPD 접근방식을  
통해 기존의 공간 소모가 많은 기판 상의 라우팅 패턴을 제  
거하고, 전반적인 전력밀도와 전력 시스템 효율을 달성할  
수 있다. WSE 외에도 이 VPD는 긴밀하게 클러스터링 된  
프로세서 구성에서 더 높은 전력을 분배할 수 있는 새로운  
가능성을 보여주고 있다.

## 랙 쿨링의 새로운 동향

물론 확장되는 컴퓨팅 인프라에서 점점 더 높은 컴퓨팅  
성능을 요구하는 것은 새로운 것이 아니다. 그러나 전통적  
으로 클라우드 데이터센터 분야에 주로 사용되는 프로세  
서는 공기 쿨링 기술로 서버 랙을 경제적으로 쿨링하기 위

해 전력 엔벨  
로프(Power  
Envelope)를  
제한하도록  
설계되어 왔  
다. 여러 측면  
에서 경제적  
인 서버 생산  
을 제약하는



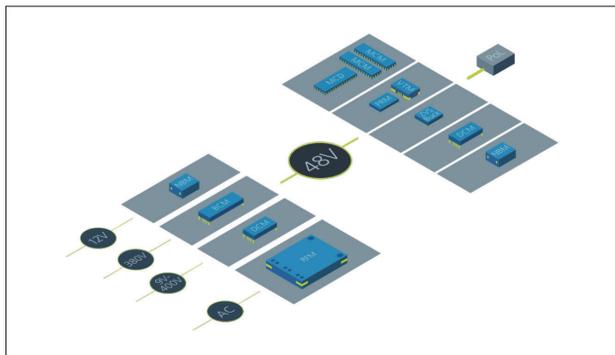
주요 요소는 프로세서 전력을 관리가 용이한 대략 200와트 이하의 임계 값으로 제한하는데 집중해 왔기 때문이다. 그러나 AI가 등장하면서 HPC는 물론, 클라우드 데이터센터에도 첨단 액체 냉각 및 침지 냉각(Immersion Cooling) 도입이 증가하고 있다.

보드 전반에 걸쳐 전력 프로세싱이 높아진 것은 2019년 오픈 하드웨어 컴퓨팅 가속기 사양인 OAM(OCP Accelerator Module)이 출시되면서 본격화된 것으로 볼 수 있다. 이 사양은 인텔과 AMD가 주도하는 업계 협력을 통해 OCP에 기여하고 있으며, 페이스북(Facebook), 마이크로소프트(Microsoft), 바이두(Baidu)와 같은 다른 업계 선두주자들의 상당한 지원을 받고 있다. OAM 사양은 독점적인 AI 하드웨어 시스템과 관련된 구현 과제와 설계 복잡성을 줄임으로써 새로운 AI 가속기 채택을 가속화하는데 목표를 두고 있다. 이미 여러 클라우드 컴퓨팅 제공업체들은 이러한 AI급 OAM 카드를 데이터센터에 구축하고 있다.

OAM이 도입되면서 데이터센터의 전용 AI 랙에 대한 새로운 아이디어를 모색하기 위해 업계에서는 많은 논의들이 전개되고 있으며, 페이스북은 이와 관련된 계획을 공개하기도 했다. 이러한 분야의 많은 공급업체들은 클라우드 데이터센터를 위한 침지 냉각 방식의 AI 랙을 면밀히 검토하고 있다. 이는 2년 전만해도 상상할 수 없었던 아이디어다.

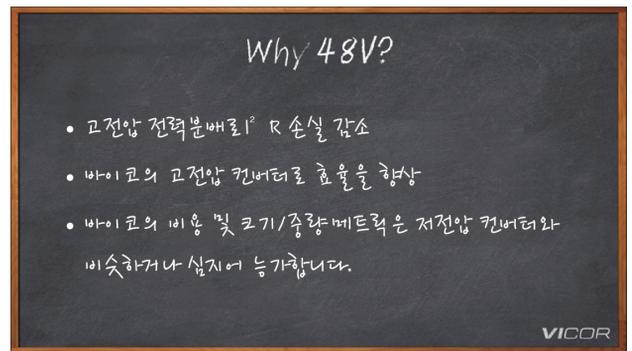
### 48V로의 진화

OAM의 또 다른 흥미로운 측면은 12V와 48V를 모두 지원할 수 있도록 설계되었다는 점이다. 따라서 기존의 12V 데이터센터 인프라의 지속적인 요구를 수용하면서도 미래



지향적인 48V 채택 요구까지도 지원할 수 있다. 그러나 대부분의 OAM 고객들은 12V에 비해 추가적으로 전력 향상을 달성할 수 있는 48V 설계에 주력할 것으로 예상된다.

이러한 고객들은 구글의 지원으로 이전에 OCP에서 도입한 48V 서버 및 분산형 인프라 표준을 부분적으로 따르고 있다. 기존의 12V 아키텍처와 비교해 48V 버스 아키텍처는 시스템 엔지니어들이 보다 높은 변환 효율과 높은 전력밀도 및 낮은 분배 손실을 제공하는 시스템을 구현할 수 있도록 해준다. 따라서 매크로 레벨에서 48V 데이터센터 서버 인프라는 에너지 손실을 30% 이상 줄일 수 있기 때문에 클라우드 데이터센터 제공업체들은 12V에서 48V로 전환하기 위한 노력을 가속화하고 있다.



한편 12V 인프라를 활용하는 클라우드 제공업체들은 12V에서 48V로 변환하여 98% 이상의 최대 효율을 제공하는 비절연 스텝업 컨버터(Step up Converter) 옵션을 이용할 수도 있다. 이를 통해 48V 전력분배로 전환하면서 차세대 AI 카드를 활용할 수 있는 유연성을 확보할 수 있다. 또한 최근 등장하고 있는 양방향 48V/12V 컨버터는 클라우드 데이터센터 제공업체들이 진화하는 인프라에 따라 두 표준을 모두 지원하거나 또는 하나만 지원할 수 있는 부가적인 유연성을 제공한다.

이러한 동향과 많은 변화들은 2020년으로 접어들면서 AI와 수퍼컴퓨팅, 클라우드 데이터센터 커뮤니티 사이에서 최우선 과제가 되고 있다. 이전에는 완전히 별개의 영역으로 간주되었지만, 최근에는 전력과 쿨링 기술이 이전에 생각했던 것보다 더 많은 공통점을 가지고 있다는 것을 확인할 수 있다. **SNV**