**Hewlett Packard Enterprise**

# Deep Learning at HPE

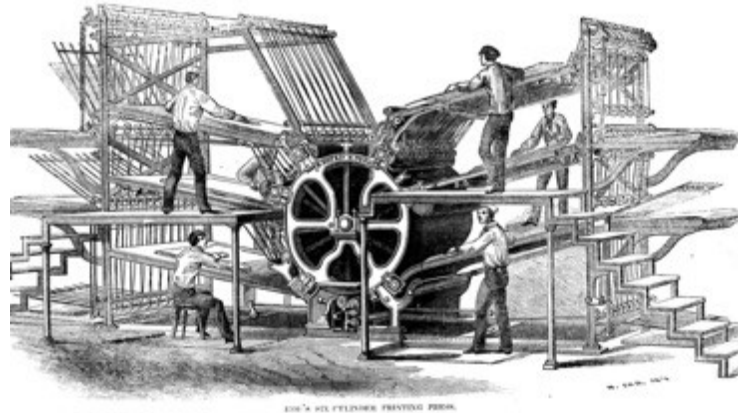# Are we on the brink of a ....

**Change 1:**
Moving from gather and hunting to settling down to farms and ports



**Change 2:**
Developing the printing press and industrial revolution



**Latest Change:**
The greatest change of our lives. Artificial Intelligence



**Hewlett Packard**
Enterprise

# Where would the road take us?

Advances in artificial intelligence will transform modern life by reshaping transportation, health, science, finance, and the military.

"High-level machine intelligence" (HLMI) is achieved when unaided machines can ac- complish every task better and more cheaply than human workers.
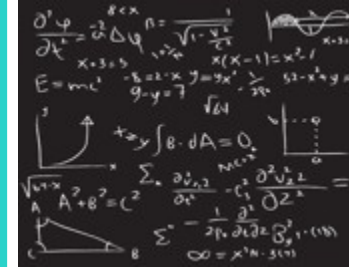
Driving a truck - 2027

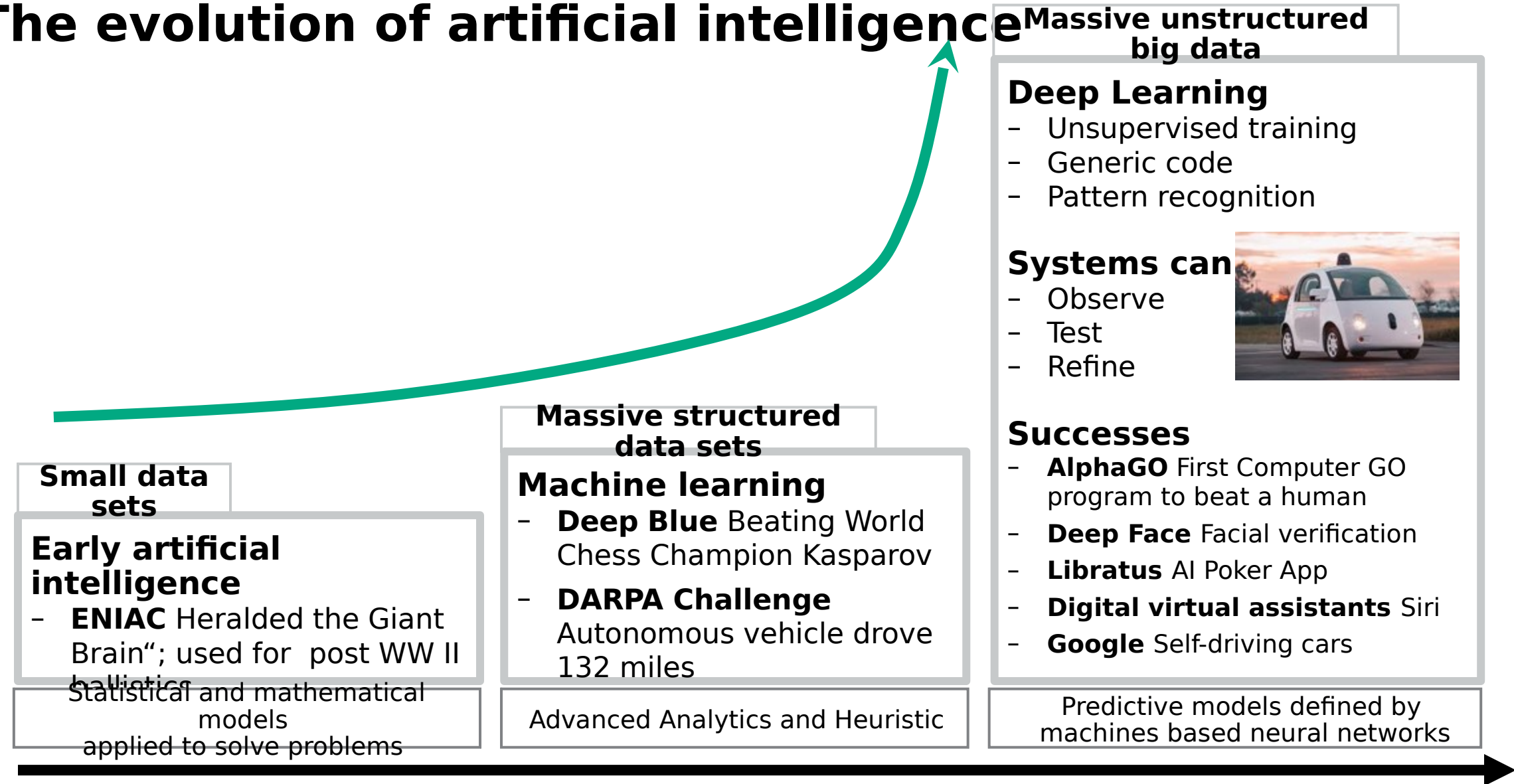Retail - 2031

Surgeon - 2043

Writing a bestseller – 2049

Math Research - 2060

Full Automation of labor – 2140

Grace et al , When Will AI Exceed Human Performance? Evidence from AI Experts

# The evolution of artificial intelligence



**Massive unstructured big data**

**Deep Learning**
- Unsupervised training
- Generic code
- Pattern recognition

**Systems can**
- Observe
- Test
- Refine

**Successes**
- **AlphaGO** First Computer GO program to beat a human
- **Deep Face** Facial verification
- **Libratus** AI Poker App
- **Digital virtual assistants** Siri
- **Google** Self-driving cars

**Massive structured data sets**

**Machine learning**
- **Deep Blue** Beating World Chess Champion Kasparov
- **DARPA Challenge** Autonomous vehicle drove 132 miles

**Small data sets**

**Early artificial intelligence**
- **ENIAC** Heralded the Giant Brain"; used for  post WW II ballistics

Statistical and mathematical models
applied to solve problems

Advanced Analytics and Heuristic

Predictive models defined by machines based neural networks

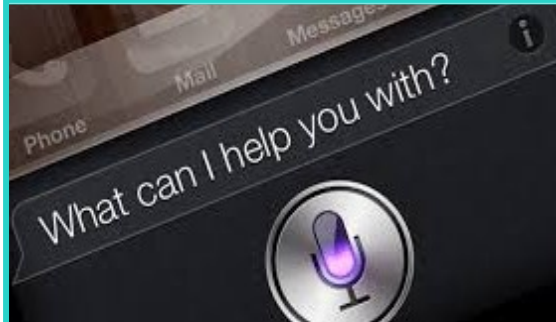**1940 – 1980**

**1990 – 2000s**

**Today**
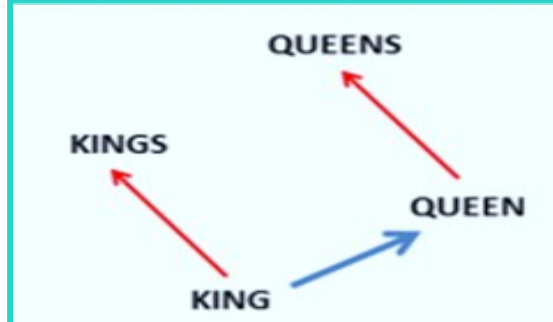
# Why deep learning?
## Applications



**Vision**

- Search & information extraction
- Security/Video surveillance
- Self-driving cars
- Medical imaging
- Robotics



**Speech**

- Interactive voice response (IVR) systems
- Voice interfaces (Mobile, Cars, Gaming, Home)
- Security (speaker identification)
- Health care
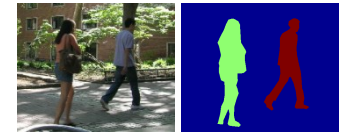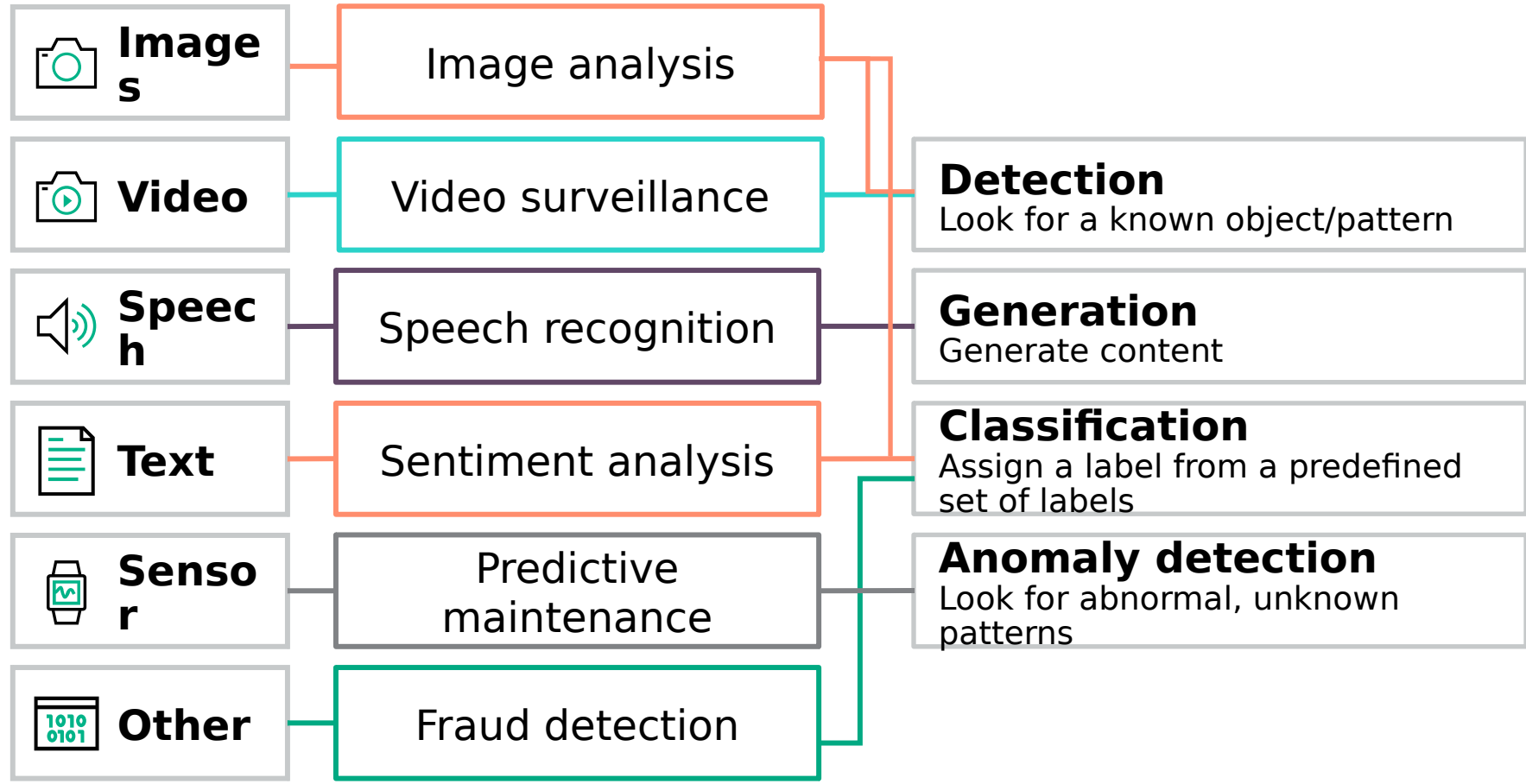- People with disabilities



**Text**

- Search and ranking
- Sentiment analysis
- Machine translation
- Question answering



**Other**

- Recommendation engines
- Advertising
- Fraud detection
- AI challenges
- Drug discovery
- Sensor data analysis
- Diagnostic support

# Applications break down

| | | |
|---|---|---|
| **Images** | Image analysis | **Detection**<br>Look for a known object/pattern |
| **Video** | Video surveillance | |
| **Speech** | Speech recognition | **Generation**<br>Generate content |
| **Text** | Sentiment analysis | **Classification**<br>Assign a label from a predefined set of labels |
| **Sensor** | Predictive maintenance | **Anomaly detection**<br>Look for abnormal, unknown patterns |
| **Other** | Fraud detection | |

# How an individual customer's AI evolves

## Explore
### How can AI help me?

**Do things better**
- Product development
- Customer experience
- Productivity
- Employee experience

**Do new things**
- New disruptions

## Experiment
### How can I get started?

**Boundary** constraints (regulations, etc.)

**Data**

Data model? Location?

How to **create** a model?
- Homegrown solution or open source?
- Simple ML or scalable DL?

**Design**

How to design and deploy the PoC?
- On-prem, cloud?
- How to think about inference

**Performance**

What is the best config to run?
How to tune the model to improve accuracy?

## Scale up and Optimize
### How can I scale and optimize?

**Provisioning** for inference

**Infrastructure scale up**
- Training
- Inference
- On-prem / cloud / hybrid

**Data management**
- Between edge and core
- Security
- Updates
- Regulations
- Tracing

**Hewlett Packard Enterprise**

# What about AI consumers ?

## Do it yourself

Current wave of AI / Machine Learning is core to their business. All in-house

**Google, Baidu, Facebook, Microsoft, Apple, etc.**

## How do I do it ?

Could benefit from better data science, machine learning, but it is not historically their core-competency

**Banks, advertisers, healthcare, manufacturing, food, automotive, etc.**

Not ready for an ASIC. Don't know what they need exactly. Many still developing on CPUs. Can't use solutions that can't be verified or understood

## I know better

Super-Experts – current wave is woefully inadequate

**Government – DoD, DoE, NSA, NASA, etc.**

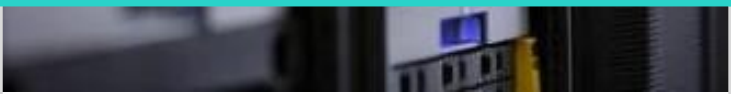Begging for higher performance ASICs. Know exactly what they want to do. Strong technology pull.

# Key IT challenges are constraining deep learning adoption
Limited knowledge, resources and capabilities

| How to get started? | How to go to production? | How to scale and optimize? |
|---|---|---|
| **Introducing the Deep Learning Cookbook** | | **HPE - Novumind Improving Deep Learning Scalability** |
| *"I need simple, infrastructure and software capabilities to rapidly and efficiently support deep learning app development."* | *"I could use more expert advice and tailored solutions for migrating and integrating apps in a production environment."* | *"I need help integrating the latest technologies into my deep learning environment to accelerate actionable insights."* |
| **Immature, sub-optimal foundation** | **Inability to scale and integrate** | **Lack of technology integration capabilities** |

# HPE and Novumind

Hewlett Packard
Enterprise
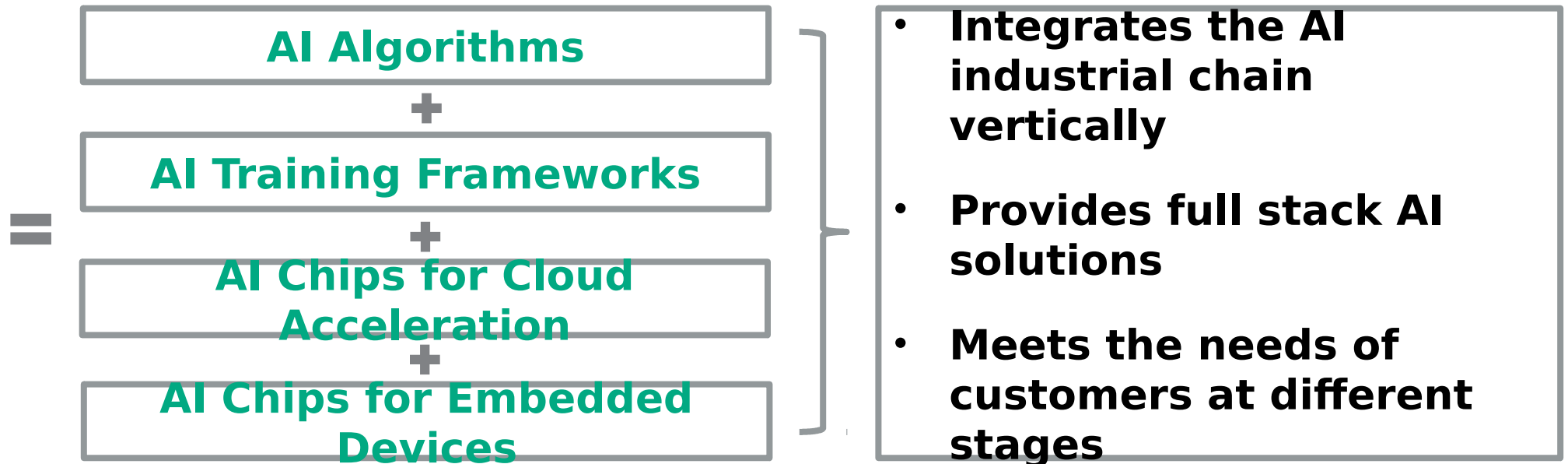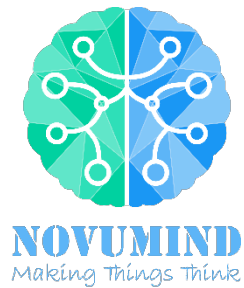
# NovuMind
## – Providing full stack AI solutions

NovuMind is a multinational AI technology company, headquartered in the heart of Silicon Valley, with branch offices in Beijing, Hong Kong, Guangzhou and Taipei.

**NOVUMIND**
*Making Things Think*

**=**

| AI Algorithms |
|:---:|

**+**

| AI Training Frameworks |
|:---:|

**+**

| AI Chips for Cloud Acceleration |
|:---:|

**+**

| AI Chips for Embedded Devices |
|:---:|

- **Integrates the AI industrial chain vertically**

- **Provides full stack AI solutions**

- **Meets the needs of customers at different stages**

# Which problems are we trying to solve ?

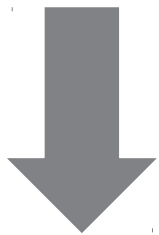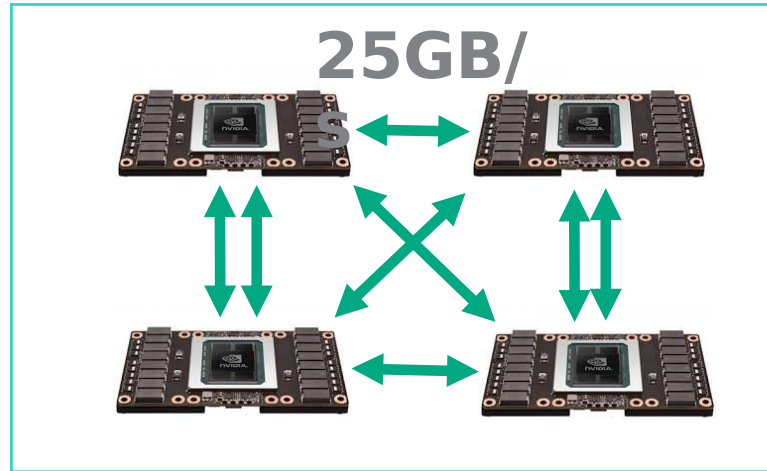| | |
|---|---|
| Need to adapt system to growing needs & data | • **FLEXIBLE** |
| Processing data, selecting framework & network and launching a job on any number of GPU must be easy | • **EASY** |
| Deep learning must run at optimal speed when system evolves / Hardware resources must be fully utilized | • **EFFICIENT** |
| Adding server must not degrade performances | • **SCALABLE** |

**Hewlett Packard**
Enterprise

# Flexibility: Topology does matter
Looking at bandwidth inter & intra node

**16GB/s**

**25GB/s**

**16GB/s**

**4:1**

**10GB/s**

**10GB/s**

**10GB/s**

**10GB/s**

**16GB/s**

**8:1**

**10GB/s**

**10GB/s**

# GPU RDMA : efficient communication to scale out



**OSU benchmark / Device to Device / Bandwidth**

Bandwidth (MS/s)

message size (KB)

Almost saturating
the interconnect

Intel OPA    MLX EDR

# Easy: Job management Web User Interface

# Easy: Resource Monitoring

# Efficient

- Latest Nvidia GPU already supported on HPE servers

- CuDNN
- NCCL

Hardware

Software

Start-of-art techniques

Tuning

- GPU RDMA

- Numa optimisation
- CPU / GPU pinning

# Efficient: Leveraging Novumind experience for optimal runs

- Data Augmentation
  - scale and aspect
  - Color
  - Weight decay
- Per epoch data shuffling
- Base on Novumind's domain experience of predefined set of meta-parameter, tuning become cook book recipes.
- In certain scenario, LR scarified a little bit accuracy for much faster convergency
- Expert knowledge of past experience to tune neural networks. For example, in Image classification, No need to search for potential tuning. Novuforce will suggest optimal ones.
- Different verticals optimal parameters ( security, heathcare, transportation, financial services)

**Hewlett Packard**
Enterprise

# Scalability: Benchmark Results



NovuCaffe Benchmark (fix batch size per GPU)

| Iterations | Model | Batch Size | GPUs | Elapsed (s) | Speedups | Final Loss |
|---|---|---|---|---|---|---|
| 2900 | resnet101 | 1 * 128 | 1 | 12027 | - | **4.6682** |
| 2900 | resnet101 | 4 * 128 | 4 | 13378 | **3.5961** | **3.4406** |
| 2900 | resnet101 | 8 * 128 | 8 | 13748 | **6.9985** | **3.0311** |
| 2900 | resnet101 | 16 * 128 | 16 | 13889 | **13.855** | **2.3069** |
| 2900 | resnet101 | 32 * 128 | 32 | 13954 | **27.5809** | **1.7068** |

Hewlett Packard
Enterprise

# Scalability: Benchmark Results



Linear Scaling of NovuForce Deep Learning

Images per Second — Number of GPUs

GoogleNet
ResNet50

# Deep Learning Cookbook

Hewlett Packard
Enterprise

# Where to start ?

## Recommend DL stack by vertical application

| | | | | | |
|---|---|---|---|---|---|
| **Verticals** | Voice interfaces | Social media | Manufacturing | Oil & gas | Connected cars |
| **Data type** | Speech | Images | Video | Sensor data | |
| **Data** | Small | Moderate | Large | | |
| **Typical layers** | Convolutional | Fully-connected | Recurrent | ... ← **Neural Network sits here** | |
| **Frameworks** | TensorFlow | Caffe 2 | CNTK | Torch | ... |
| **Infrastructure** | x86 | GPUs | FPGAs | TPU ? | ... |

**Hewlett Packard Enterprise**

# Neural Network : Popular Networks

| Network | Model size (# params) | Model size (MB) | GFLOPs (forward pass) |
|---------|----------------------|-----------------|----------------------|
| AlexNet | 60,965,224 | 233 MB | 0.7 |
| GoogleNet | 6,998,552 | 27 MB | 1.6 |
| VGG-16 | 138,357,544 | 528 MB | 15.5 |
| VGG-19 | 143,667,240 | 548 MB | 19.6 |
| ResNet50 | 25,610,269 | 98 MB | 3.9 |
| ResNet101 | 44,654,608 | 170 MB | 7.6 |
| ResNet152 | 60,344,387 | 230 MB | 11.3 |

# Hardware : Today's scale and needs
## Model size, data size, compute requirements

| Application | Model | Training data | FLOPs per epoch |
|---|---|---|---|
| Vision | $1.7 * 10^9$<br>~6.8 GB | $14*10^6$ images<br>~2.5 TB (256x256)<br>~10 TB (512x512) | $6*1.7*10^9*14*10^6$<br>~$1.4*10^{17}$ |

1 epoch per hour:
~39 TFLOPS

**Today's hardware:**

NVIDIA Tesla V100: 15 TFLOPS SP (30 TFLOPS FP16 , 120 TFLOPS Tensor ops), 12 GB memory

NVIDIA Tesla P100: 10.6 TFLOPS SP, 16 GB memory

NVIDIA Tesla K40: 4.29 TFLOPS SP, 12 GB memory

NVIDIA Tesla K80: 5.6 TFLOPS SP (8.74 TFLOPS SP with GPU boost), 24 GB memory

INTEL Xeon Phi: 2.4 TFLOPS SP

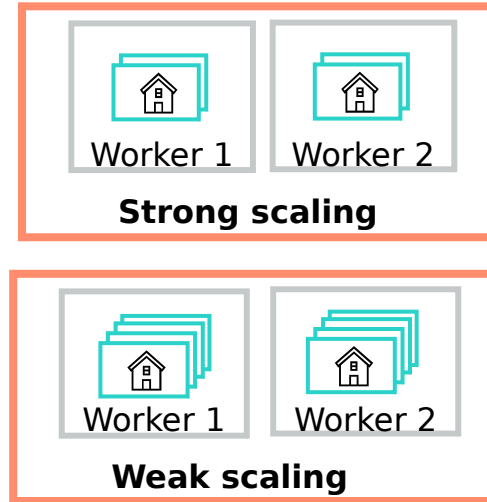Superdome X: ~21 TFLOPS SP, 24 TB memory

# So what to recommend?

# Building performance models

Alex Net

GoogleNet

VGG-16, VGG -19

ResNet 50, 101,152

Eng Acoustic Model

TensorFlow

Caffe 2

Tensor RT

BVLC Caffe

Worker 1    Worker 2

**Strong scaling**

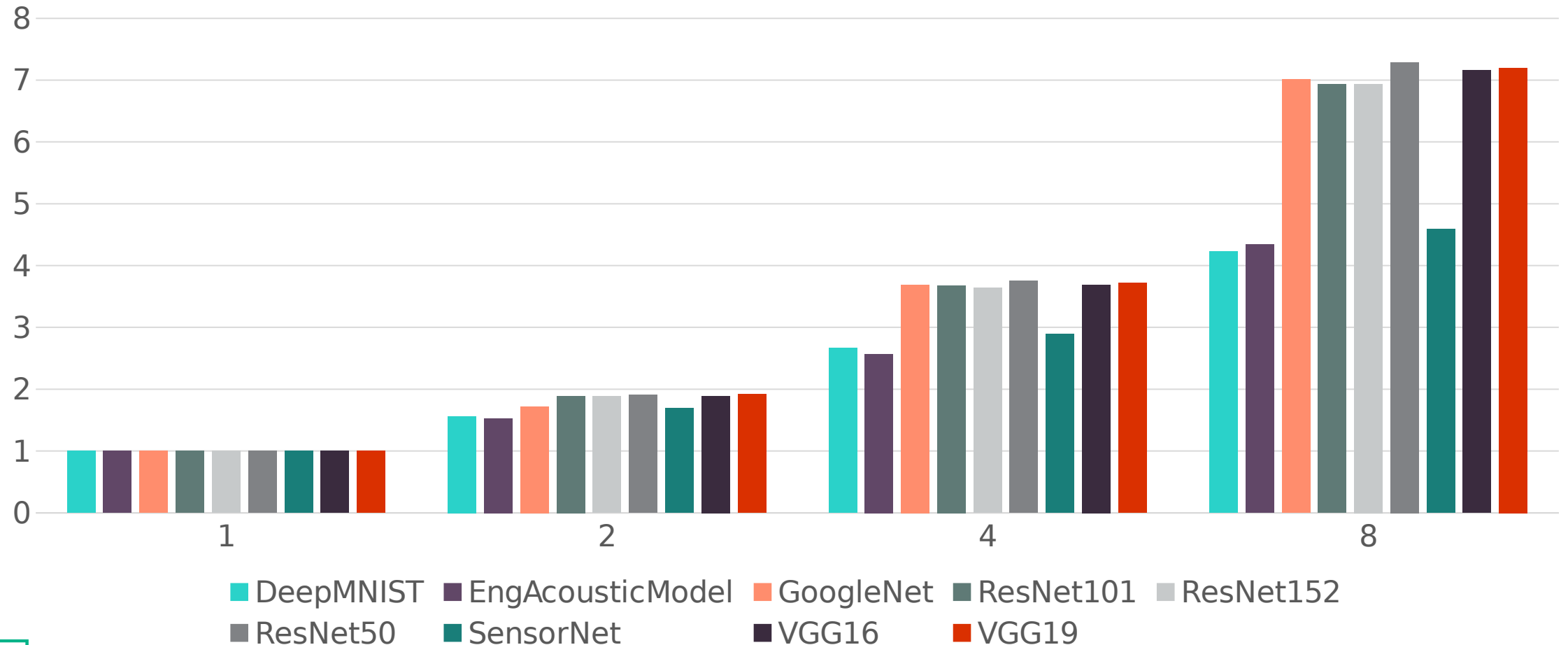Worker 1    Worker 2

**Weak scaling**

**Hardware**

Scalable, automated
real-time intelligence

Populated with 8 GPUs

# TensorFlow – Weak Scaling – Training – Different models perfromance in Tensor Flow . Scaling up to 8 GPUs
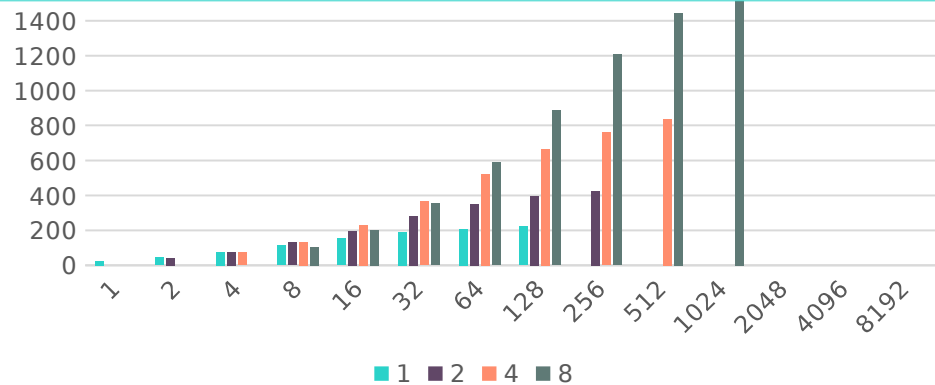
Speedup for up to 8 GPUs

# TensorFlow - Inference ( Inferences per Second) - Different Models witth different Batch numbers

DeepMNIST

GoogleNet



HOW TO ANALYZE ALL THE DIFFERENT NUMBERS .

AS WE ADD MORE OPTIONS and MORE TECHNOLOGIES IT WOULD BE IMPOSSIBLE TO USE



**Hewlett Packard Enterprise**

# Select ideal technology configurations
## with HPE Deep Learning Cookbook

### "Book of recipes" for deep learning workloads



- **Comprehensive tool set** based on extensive benchmarking
- **Includes** various models with 8 DL frameworks and 8 HPE hardware systems
- **Estimates workload performance** and recommends an optimal HW/SW stack for that workload

### Expert advice to get you started



- **Informed decision making** - optimal hardware and software configurations
- **Eliminates the "guesswork"** - validated methodology and data
- **Improves efficiency** - detects bottlenecks in deep learning workloads
- **Determine scalability**

### Availability of complete toolset



- **Deep Learning Benchmarking Suite:** available on GitHub  Dec 2018
- **Deep Learning Performance Analysis Tool:** planned to be released in the beginning of 2018**.**
- **Reference configurations:** available soon on HPE.com website

# Deep Learning Cookbook

**Automatic Meeting Notes**  **Video Surveillance**  **Hospital Smart Care Unit**  **Custom**

- ☐ Images
- ☑ Videos
- ☐ Text
- ☐ Speech
- ☐ Sensor Data

- ☐ Classification
- ☑ Detection
- ☐ Generation
- ☐ Anomaly Detection

- ◉ Training
  - ○ Large
  - ◉ Medium
  - ○ Small
- ○ Inference

**Recommend**

## Data and Model

| Data size | Epochs | Model |
|---|---|---|
| 10000000 | 50 | VGG19 ▽ |

## Hardware

| Server | Processor unit |
|---|---|
| Apollo 6500 ▽ | NVIDIA P100 ▽ |

| Count | Cluster size | Interconnect |
|---|---|---|
| 8 ▽ | 2 | InfiniBand FDR ▽ |

## Software

| Framework | Batch size | Scaling |
|---|---|---|
| Caffe2 ▽ | 1024 ▽ | strong ▽ |

**Add**

## Training performance



training time vs devices (1–16)

| | Data | | | Hardware | | Software | Time (hours) | |
|---|---|---|---|---|---|---|---|---|
| 🟧 | Size 10000000 | Epochs 50 | Model AlexNet | Server Apollo 6500 / Count 8 | PU NVIDIA P100 / Cluster size 2 / Interconnect IB | Framework Caffe2 / Batch 1024(strong) | 22.5 | ✕ |
| 🟪 | Size 10000000 | Epochs 50 | Model GoogleNet | Server Apollo 6500 / Count 8 | PU NVIDIA P100 / Cluster size 2 / Interconnect IB | Framework Caffe2 / Batch 1024(strong) | 26.3 | ✕ |
| 🟩 | Size 10000000 | Epochs 50 | Model VGG19 | Server Apollo 6500 / Count 8 | PU NVIDIA P100 / Cluster size 2 / Interconnect IB | Framework Caffe2 / Batch 1024(strong) | 147.4 | ✕ |

**Remove all**

Hewlett Packard
Enterprise

Hewlett Packard
Labs

# Thank you

**Natalia Vassilieva**
nvassilieva@hpe.com

**Sergey Serebryakov**
sergey.serebryakov@hpe.com

**Sorin Cheran**
sorin.cheran@hpe.com

**Bruno Monnet**
bruno.monnet@hpe.com

**Hewlett Packard**
Enterprise