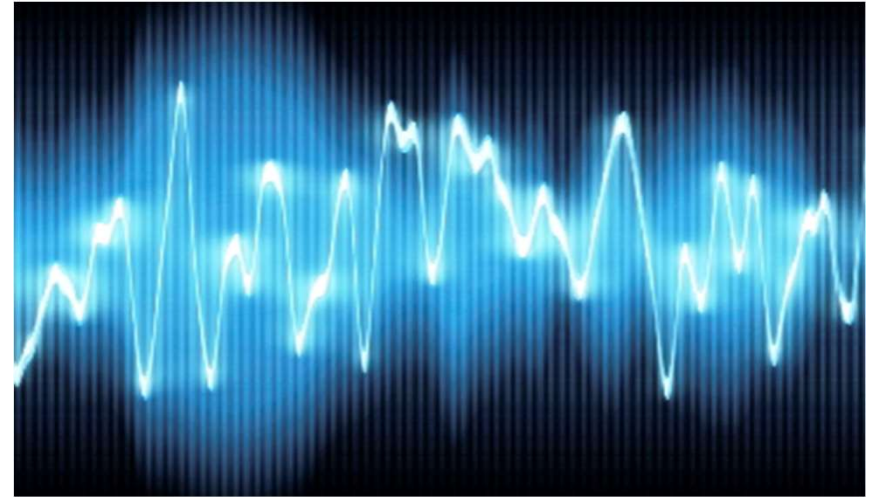# Deep Learning Revolution: From Theory to Impact





## Chris Rowen

CEO
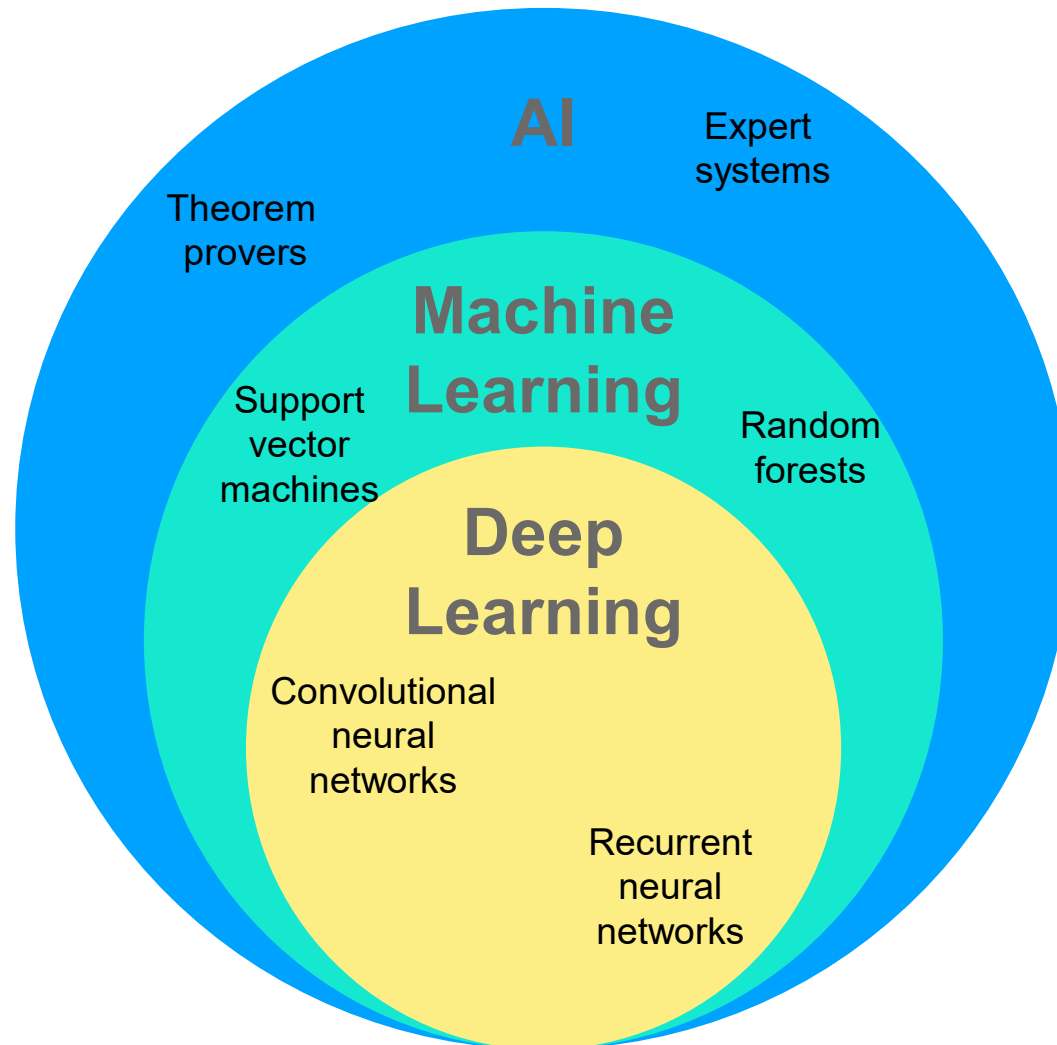Babblelabs Inc.
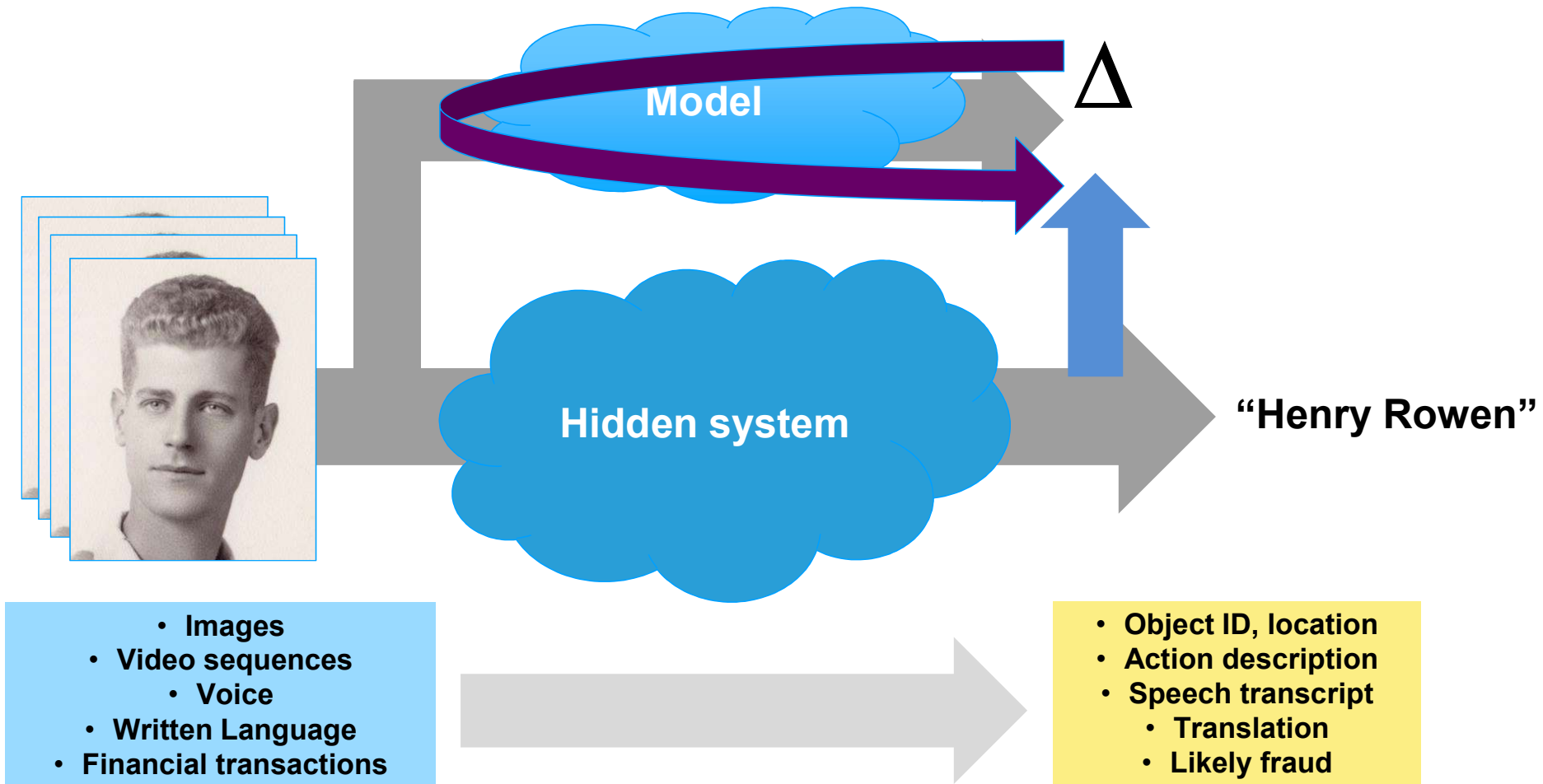
EDPS

September 2018

## Hype Metrics:

- One Google page hit on "AI" for every person in US + India + China
- 11,300 "artificial intelligence" startups [CrunchBase]
- 16,500 papers on "neural network" on arxiv.org – most in past 24 months

babblelabs

# Quick Taxonomy



AI

Expert systems

Theorem provers

**Machine Learning**

Support vector machines

Random forests

**Deep Learning**

Convolutional neural networks

Recurrent neural networks

babblelabs

# Deep Learning Foundations

*The construction of a complex numerical model that mimics the behavior of a very complex but hidden system:*

Model

Hidden system

$\Delta$

"Henry Rowen"

- Images
- Video sequences
- Voice
- Written Language
- Financial transactions

- Object ID, location
- Action description
- Speech transcript
- Translation
- Likely fraud

babblelabs

# Vision is Fundamentally Hard

- Big computation in embedded inference, **huge** in training (but less frequent)

- Typically need large *labeled* data-sets

- Example: ImageNet Classification:

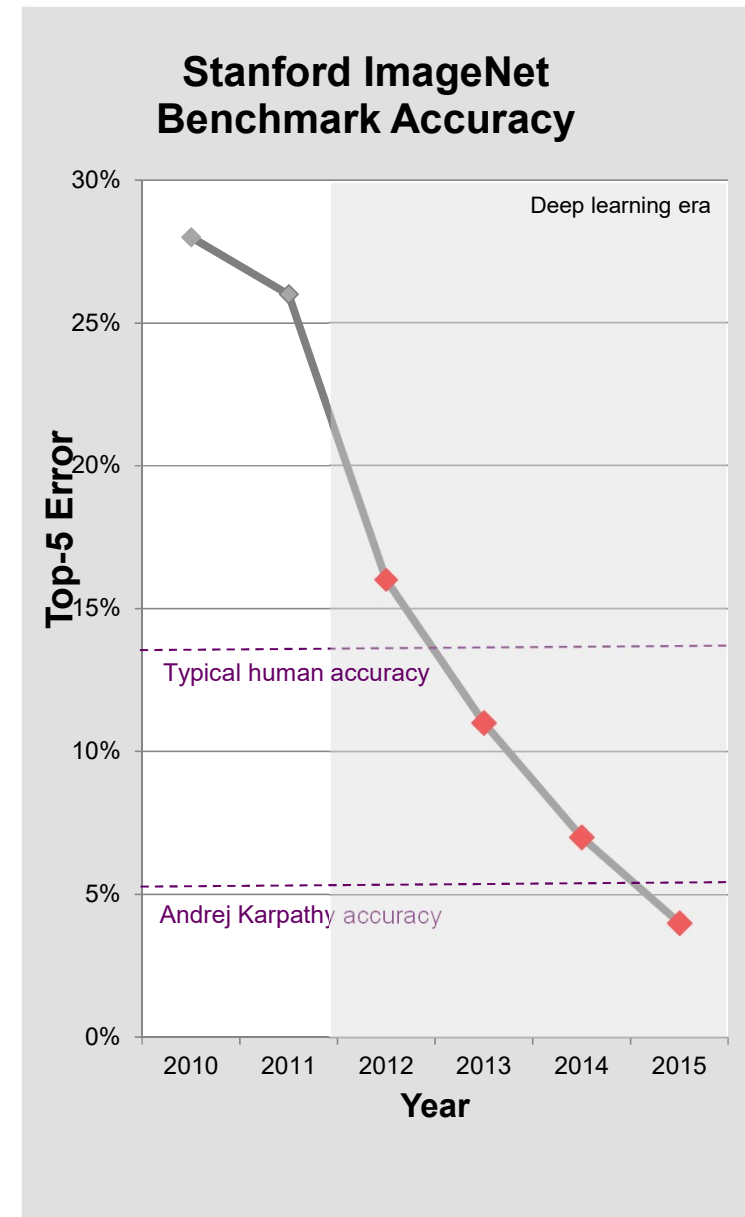  - 1.2M images

  - 1000 categories

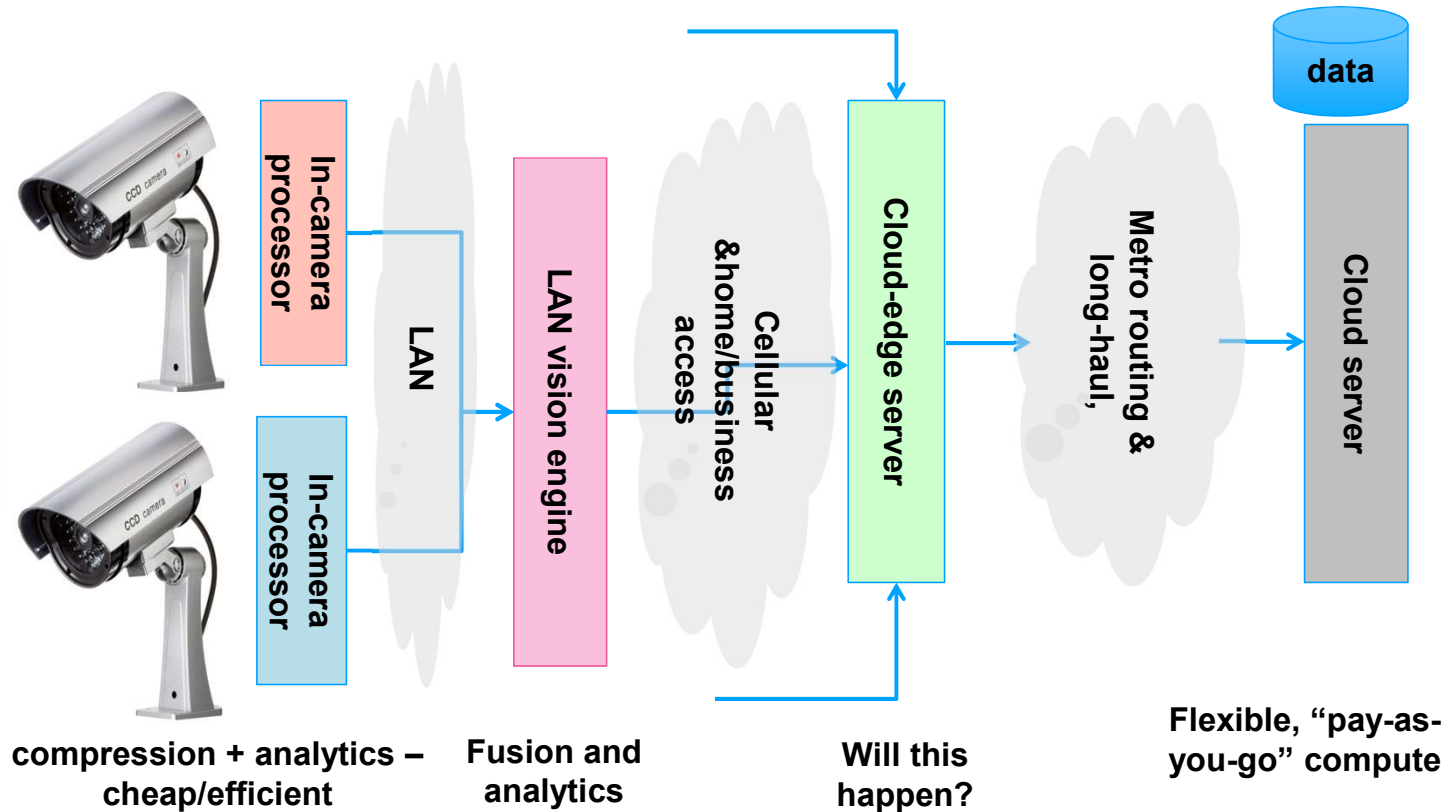  - 120 breeds of dogs

**Tibetan mastiff**     **Shih-Tzu**     **Norwegian elkhound**

**Stanford ImageNet Benchmark Accuracy**

Deep learning era

Typical human accuracy

Andrej Karpathy accuracy

Top-5 Error

Year

# Where will we put the "smarts"?



**In-camera processor**

**In-camera processor**

**LAN**

**LAN vision engine**

**Cellular &home/business access**

**Cloud-edge server**

**Metro routing & long-haul,**

**Cloud server**

**data**

**compression + analytics – cheap/efficient**

**Fusion and analytics**

**Will this happen?**

**Flexible, "pay-as-you-go" compute**

System responsiveness

Scope of data analysis

Privacy

Computing and network cost

babblelabs

# Security, Robustness and Privacy

- More risks for privacy: more cameras, more correlation to other streams
- Mission-critical surveillance susceptible to new attacks

The usual network and device attacks

+ Database manipulation to inject bias
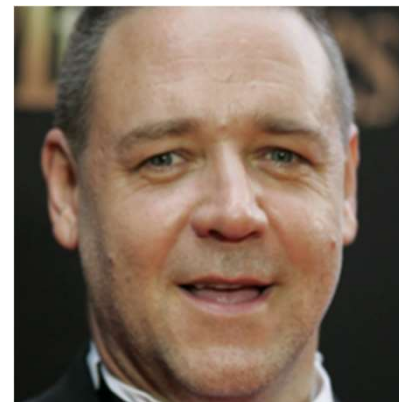
+ Classification spoofing

- Example:  Spoofing facial recognition

**Reese Witherspoon**

**Reese Witherspoon in patterned frames**
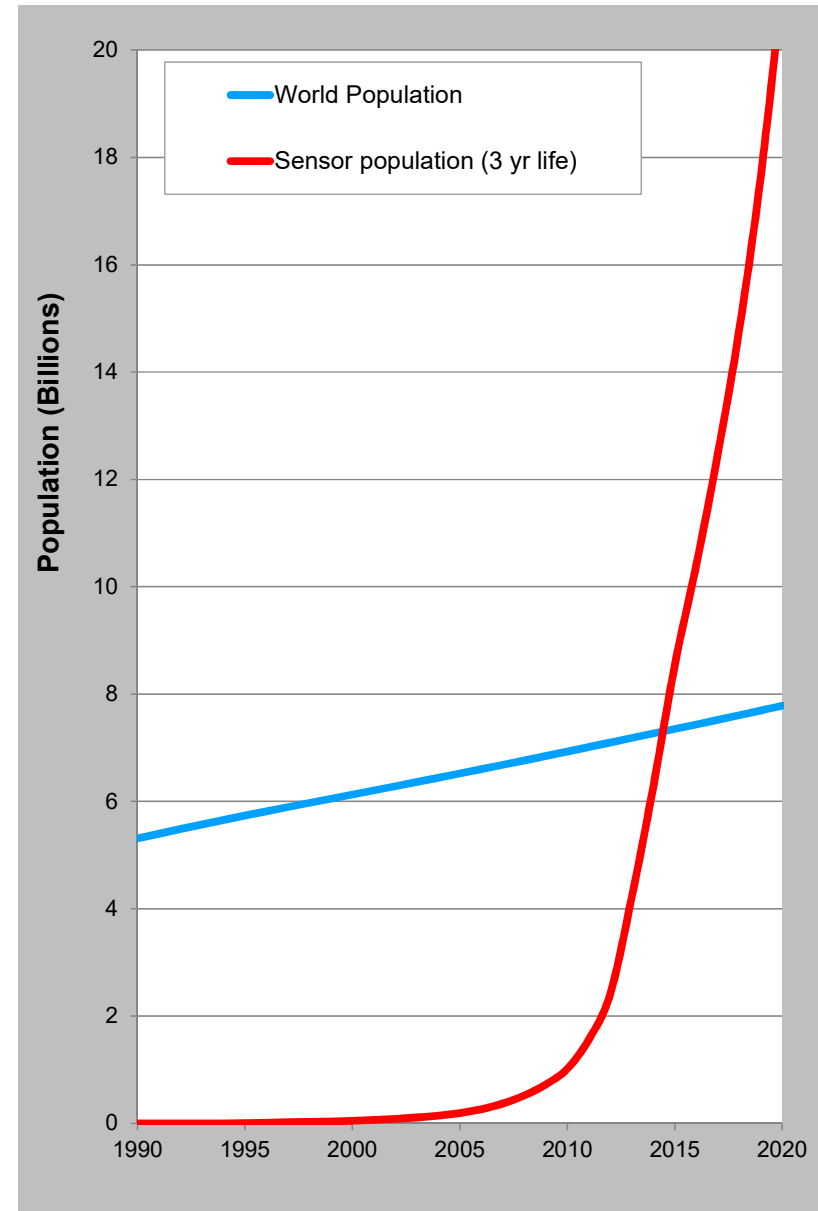
**Recognized as Russell Crowe**

**Sharif et al: Accessorize to a Crime: Real and Stealthy Attacks on State-of-the-Art Face Recognition [CMU]**

# Vision:
## *the pixel explosion*

- Rapid replacement of traditional vision by deep learning

- From 2015: more image sensors than people

- 99% of all new raw data is pixels (rest is audio)

- Massive bandwidth:

  $2\times10^{10}$ sensors * $5\times10^{8}$ pixels/sec = $10^{19}$ raw pixels/s

ZOSI 1/3" 1000TVL 960H Security Surveillance CCTV HD Camera Had IR Cut 3.6mm Lens Outdoor Weatherproof Day Night Vision 65ft Distance White
by ZOSI
$13⁹⁹ ✓prime
In Stock

★★★★☆ ▾ 247
**Product Features**
... Please noted: this *camera* did not come with power supply &video power ...

$**13⁹⁹** ✓prime
In Stock



babblelabs

SemiCo 2014

# How much does a $10 camera cost?

| camera ➔ cloud H.264 compress | Camera | Power | Network | Storage | Compute | 3 year total cost per camera |
|---|---|---|---|---|---|---|
| | HD camera w/compression | At 1$/watt-year | DSL/cable network @$10/TB | Rolling 1 day of data @ $25/TB/month | YOLO2 object detection on AWS G3 @ $0.60/hr | |
| @ 60 fps | $10 | $9 | $4,700 | $400 | $3,300 | ~$8400 |
| @ 1fps | $10 | $9 | $80 | $7 | $55 | ~$165 |
| @ 0.1 fps | $10 | $9 | $8 | $1 | $5 | ~$35 |

**Observations:**

- *Cloud-based real-time vision requires "semantic compression" at the edge*
- *Completely autonomous analysis and action is also the lowest cost*
- *Vision at the edge is biggest deep learning silicon market*
  - *Autonomous vehicles and robots*
  - Video monitoring
  - UI and social media with AR/VR

babblelabs

# Voice is Vision
## *Speech is the most human of interfaces*
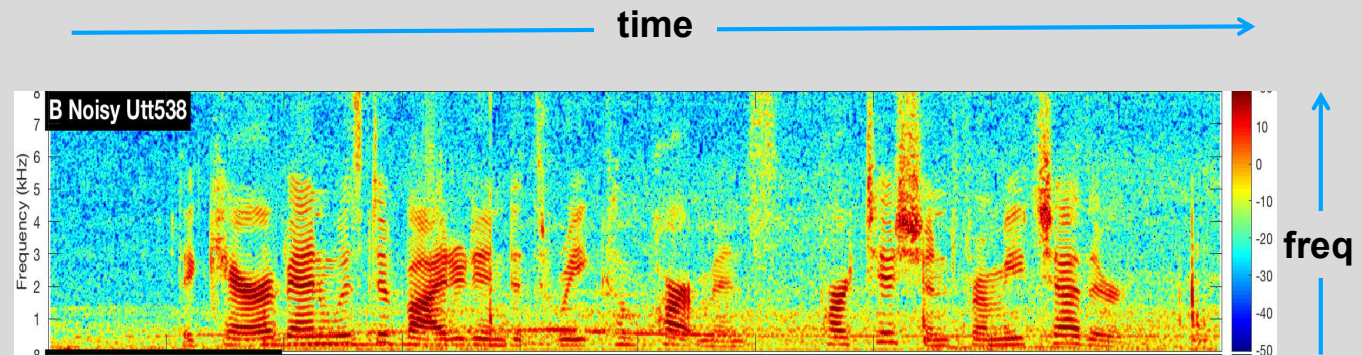
5B active electronic speech users

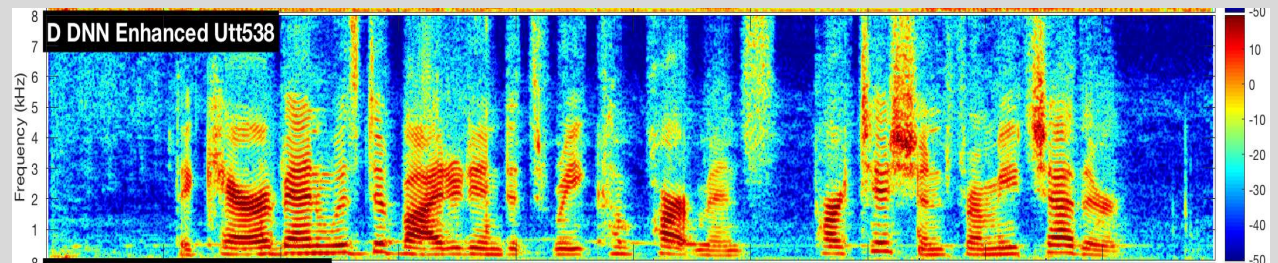>20B microphones installed by 2020

Key services:

- Noise reduction
- Speech recognition
- Speaker authentication
- Speech synthesis

**Shift to local compute**

**Transform speech signal into image: "spectrogram"**

time

B Noisy Utt538

Frequency (kHz)

freq

Vision-style neural network

D DNN Enhanced Utt538

Frequency (kHz)

babblelabs

# BabbleLabs:  Deep learning meets speech

*fresh problems, more sophisticated models, more data, more training*



0dB

**Massive Compute**

**88 TFLOPS Per Engineer**

**Massive Data Corpus**

**30,000 hr speech**
**7,000 hr music**
**3,000 hr noise**
**10,000 room models**

21dB

**Speech Enhancement**

**Speaker-specific Services**

**Speech UI Dialog**

babblelabs

# Speech Enhancement Example

**Original Video**



**Enhanced Video**



babblelabs

# Deep Learning Silicon Is Easy
*especially for inference*

$$g_j = \sigma(\sum_{i=0}^{n} C_{ji}x_i + B_j)$$

- Compute dominated by multiply-add
- Coefficients $C_{ij}$, $B_j$ read-only, heavily reused
- Memory pattern regular, static and bounded
- Programmability lets hardware span many applications
- High-level frameworks hide architecture details

babblelabs

# Deep Learning Silicon is Hard

- Impediments in efficiency:
  - mixed convolution sizes
  - non-unit strides and short,odd vector lengths
  - difficult parallelization
  - on-the-fly data reorganization
  - sparsity in coefficients and intermediate results
- Memory bandwidth;
  - large models (10s of MB)
  - fully connected layers (1 coefficient/MAC)
  - many-to-many communication between layers
  - CPU – Neural Network Engine data sharing
  - training >> inference
- MUST have comprehensive mapping, optimization and analysis SW from frameworks to silicon
- Silicon availability may be getting ahead of deployable applications

babblelabs

# Anatomy of a Neural Network Accelerator
## *Example: Google Tensor Processing Unit (TPU)*

**Neural Network Systolic Coprocessor:**
- Pumps slices of matrix multiply through array
  - Relies on CPU for control
  - 40W (TPU) + 70W (CPU)
  - Typical batch sizes: 8-200

**DDR3**

DDR3 x 2
30 GB/s

14 GB/s

30 GB/s

Intel Haswell CPU

PCle Gen3 x 16
14 GB/s

Host Interface

DDR3 Interface

**Weight FIFO/Fetcher**

30 GB/s

Data pipeline

10 GB/s

**Unified Activation Buffer 28MB**

**Systolic Data Setup**

167 GB/s

**256 x 256 8b multiplier systolic array**

Weight pipeline

Partial sum pipeline

"CISC" Instructions
Average latency: 10 cycles

Control Sequencing

**Accumulators**

**Non-linear function**

**Norm & Pool**

167 GB/s

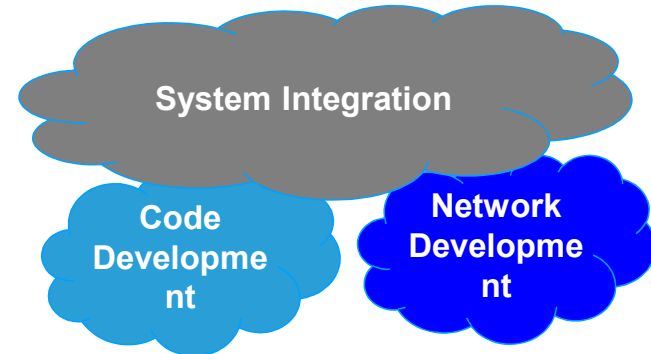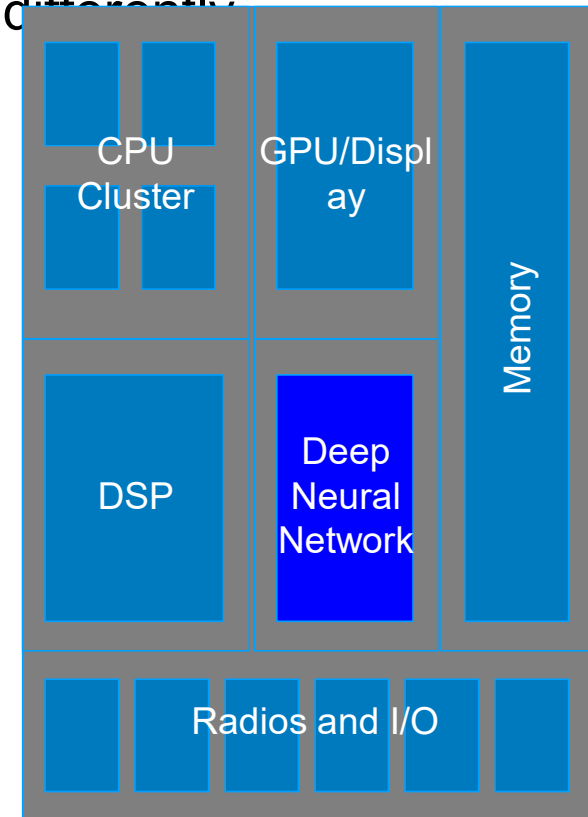# Implications for Semiconductors

1. New computing model:
   - New applications – especially vision & speech
   - new focus on data-set access
   - new development tools

2. Deep learning uses cloud and edge devices differently
   - Cloud:
     - huge compute + big memory for training
     - inference on across aggregated data
   - Edge
     - high-bandwidth real-time inference
     - minimum power and cost
     - latency- and privacy-critical applications

3. Neural networks very diverse
   - Some run on existing $\mu$controllers
   - Some CPU, GPU and DSP ➔ special compute

**System Integration**

**Code Development**

**Network Development**

CPU Cluster

GPU/Display

Memory

DSP

Deep Neural Network

Radios and I/O

babblelabs

# A Picture of Deep Learning Startups

- More than 2/3 of 372 startups focus on cloud software

- Half of all startups do vision

- Embedded dominated by vision

- Speech by startups just starting

- Many deep learning chip startups



Deep Learning AI: 372

Vision: 185

Embedded Vision: 105

Chips: 24

Embedded: 124

Speech: 26

babblelabs

Cognite Ventures

# Silicon Design Renaissance

Not Just the Big Chips and System Makers

- *Implication: high performance & low power inference will be widely available*
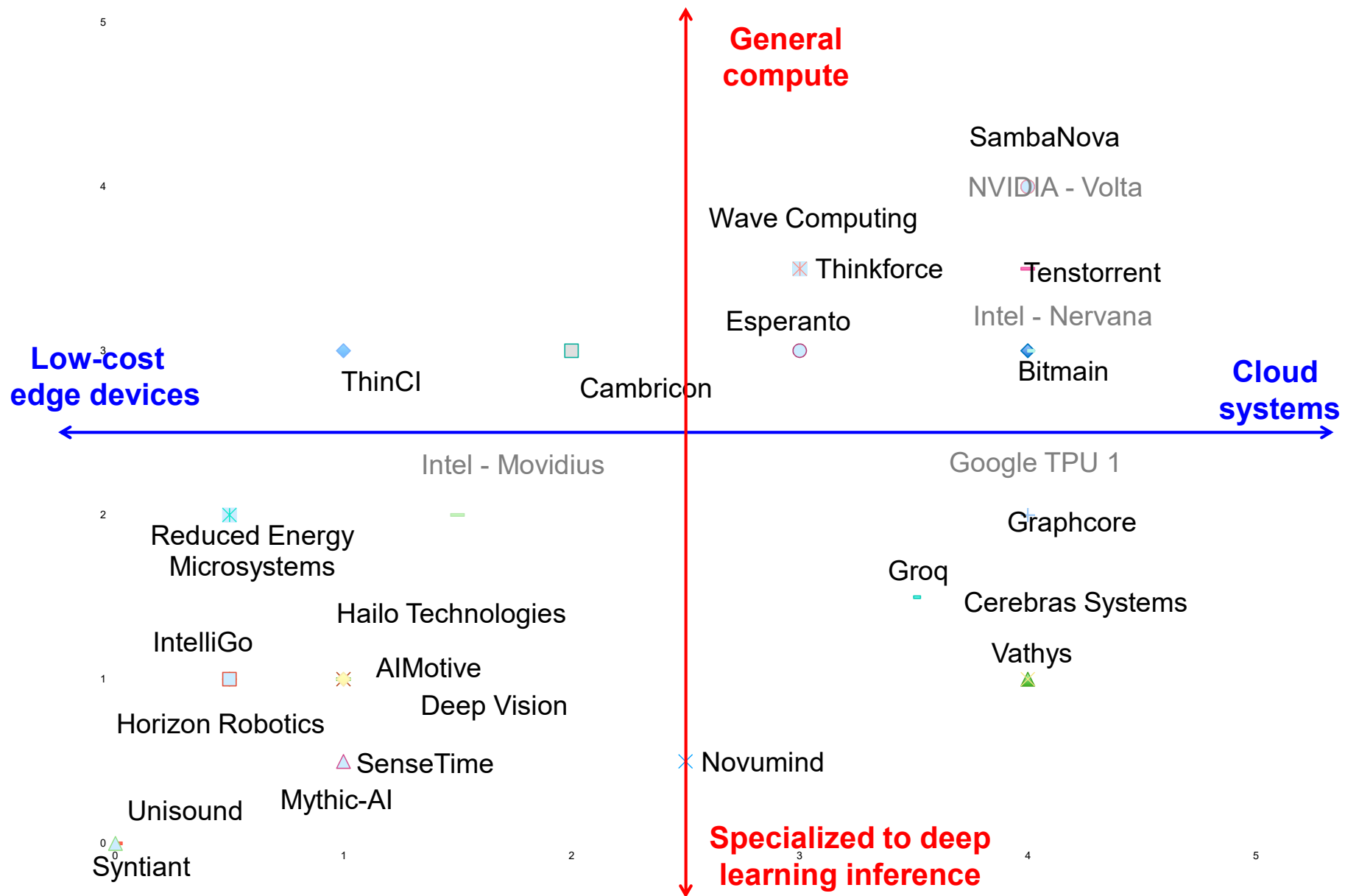
➔ In embedded devices:

- ▪ vision

- ▪ speech

➔ In cloud

- ▪ network development: training
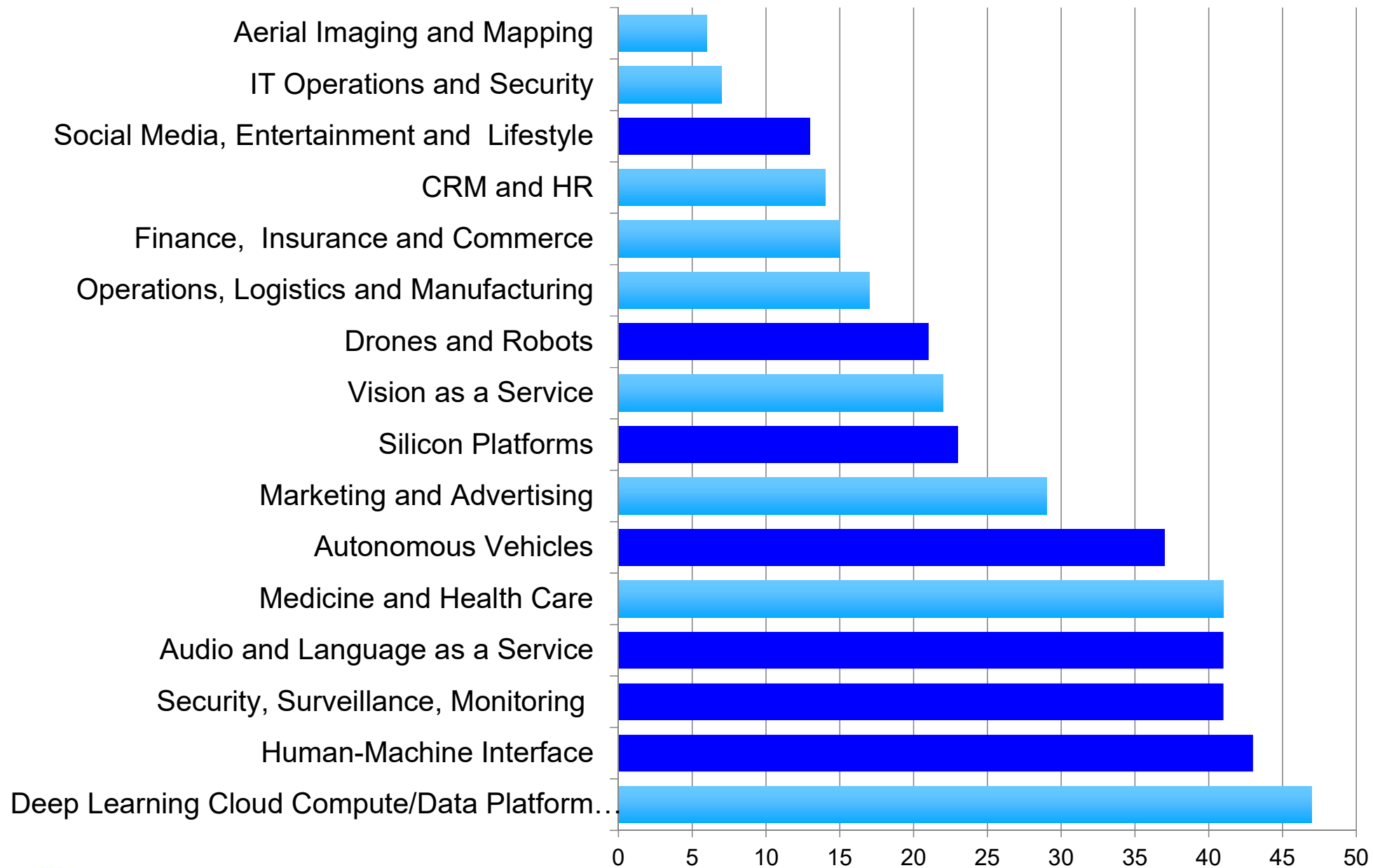
- ▪ application deployment: inference

babblelabs

| | | |
|---|---|---|
| **AIMotive** | Portable software for automated driving | Hungary |
| **Axis Semi** | Massive array of compute cores | USA |
| **Bitmain** | Coin miner builds training ASIC | China |
| **Cambricon** | Device and cloud processors for AI | China |
| **Cerebras Systems** | Specialized chip for deep-learning applications | USA |
| **Chipintelli** | Speech recognition chip for local speech processing | China |
| **Deep Vision** | Low-power silicon architecture for computer vision | USA |
| **Esperanto** | Massive array of RISC-V cores | USA |
| **FWDNXT** | Low power image recognition and classification | USA |
| **Graphcore** | Graph-oriented processors for deep learning | UK |
| **Groq** | Google spinout doing deep learning chip | USA |
| **Gyrfalcon Technology** | Low-cost, low-power, high-performance Artificial Intelligence (AI) processors. | USA |
| **Hailo Technologies** | Specialized deep learning microprocessor | Israel |
| **Horizon Robotics** | Smart Home, automotive and Public safety | China |
| **IntelliGo** | Hardware and software for image and speech processing | China |
| **Mythic-AI** | Low power NN inference IC using flash+analog+digital | USA |
| **Novumind** | AI for IoT | USA |
| **Preferred Networks** | Real time data analytics with deep learning and Chainer library | Japan |
| **Rain Neuromorphics** | Nanotechnology for AI | USA |
| **Reduced Energy Microsystems** | Lowest power silicon for deep learning and machine vision | USA |
| **SambaNova** | Coarse Grain Reconfigurable Array for matrix arithmetic | USA |
| **SenseTime** | Computer vision | China |
| **Syntient** | Customized analog neural networks | USA |
| **Tenstorrent** | Deep learning processor: designed for faster training and adaptability to future algorithms | Canada |
| **ThinCI** | vision processing chips | USA |
| **Thinkforce** | AI chips | China |
| **Unisound** | AI-based speech and text | China |
| **Vathys** | Deep learning supercomputers | USA |
| **Wave Computing** | Deep Learning computers based on custom silicon | USA |

# Sorting Out Deep Learning Silicon

# Segments for Deep Learning Startups

High potential volume

| Segment | Value |
|---|---|
| Aerial Imaging and Mapping | 6 |
| IT Operations and Security | 7 |
| Social Media, Entertainment and Lifestyle | 13 |
| CRM and HR | 14 |
| Finance, Insurance and Commerce | 15 |
| Operations, Logistics and Manufacturing | 17 |
| Drones and Robots | 21 |
| Vision as a Service | 22 |
| Silicon Platforms | 23 |
| Marketing and Advertising | 29 |
| Autonomous Vehicles | 37 |
| Medicine and Health Care | 41 |
| Audio and Language as a Service | 41 |
| Security, Surveillance, Monitoring | 41 |
| Human-Machine Interface | 43 |
| Deep Learning Cloud Compute/Data Platform… | 47 |

babblelabs

# Where Are the Deep Learning Startups?

## Deep Learning Startups by Country



Canada
5.1%

China
7.8%

Israel
11.5%

UK
16.0%

USA
44.7%

Legend:
- USA
- UK
- Israel
- China
- Canada
- Germany
- India
- France
- Japan
- Spain
- The Netherlands
- Russia
- South Korea
- Sweden
- Austria
- Switzerland
- Denmark
- South Africa
- Hungary
- Ireland
- Argentina
- Turkey
- Singapore
- Norway
- Lithuania
- Slovenia
- Belgium
- Finland
- Czech

# Understanding China's Startups

## Cloud Deep Learning Startups by Country



- USA 44.0%
- UK 19.2%
- Israel 13.2%
- Canada 4.0%
- China 4.0%

**Legend (center):**
- USA
- UK
- Israel
- Canada
- China
- Germany
- India
- France
- Spain
- The Netherlands
- South Korea
- Japan
- Switzerland
- Russia
- Austria
- Denmark
- South Africa
- Ireland
- Argentina
- Turkey
- Singapore
- Lithuania
- Slovenia
- Belgium
- Finland
- Czech

## Embedded Deep Learning Startups by Country



- USA 45.5%
- China 15.4%
- UK 10.6%
- Israel 7.3%
- Canada 7.3%

**Legend (right):**
- USA
- China
- UK
- Israel
- Canada
- Germany
- France
- Japan
- Russia
- Sweden
- India
- Spain
- Austria
- Denmark
- Hungary
- Norway
- Taiwan

babblelabs

# Deep Learning Startups
## *A window into the future of electronics*

**Cognite** Ventures

# THE COGNITE 350
## TOP STARTUPS IN DEEP LEARNING **FEBRUARY 2018**

AMERICAS 177
ASIA 39
EUROPE MIDDLE EAST AFRICA 133

## DEEP LEARNING CLOUD COMPUTE/DATA PLATFORM AND SERVICES

DataRobot, H2O.ai, Cogital, vicarious, NanoNets, Digital Reasoning, OpenAI, minds.ai, ARIMO, groq inc., Numenta, vertex.ai, deepsense.io, Cirrascale, sentient, skymind, DATALOGUE, rapidminer, ELEMENT AI, KIMERA, loop.ai, SIGOPT, naralogics, Cycorp, FLOYD, 4Paradigm, Preferred Networks, Arya.ai, hocrox, diffblue, Intellegens, AURORA-AI, nnaisense, NEURENCE, seldon, twentybn, artelnics, deepomatic, Deep Solutions, SPARKBEYOND, KUZNECH, RAINBIRD

## MEDICINE AND HEALTH CARE

imagia, AiCure, Atomwise, Cure Metrix, ZEPHYR HEALTH, Numerate, deep genomics, DEEP6 AI, RECURSION, babblabs, SENTER, GRAIL, Butterfly Network, SENTRIAN, med what, pulseData, MEDANN, LUMINIST, BAYLABS, RECURSION, CLOUDMEDX, CareSkore, enlitic, VUNO, iCarbonX, Lunit, AVALON Ai, KHEIRON, ContextVision, foodvisor, camereyes, Intellegens, BenevolentAI, zebra, MAGENTIQ, MedyMatch

## FINANCE, INSURANCE AND COMMERCE

CEREBELLUM CAPITAL, Dataminr, KENSHO, it.ai, Tick, SI, iSENTIUM, zest, DEEP, alpha-i, FEATURE SPACE, BMLL, oseven, T

## MARKETING AND ADVERTISING

Affectiva, DEEPVISION, Seamless.AI, XIX, VUE.AI, iSENTIUM, Amplero, AUTOMAT, PERSADO, PRODUCTAI, SOTERIA INTELLIGENCE, ditto, Crossing Minds, tamr, LAYER 6, netra, PRODUCTAI, COGNITIV, AYLIEN, NUDGR, VISII, viisights, Revuze, CORTEXICA, WebyClip, Voyager Labs, real eyes, phrasee

## VISION AS A SERVICE

clarifai, Crowd AI, eyeris, MOD9, Matroid, PRISMA, Sighthound, tend, Hubino, HYPERVERGE, SeetaTech, TUPUTECH, YITU, SENSETIME, Birds.ai, deepomatic, inoven, VILYNX, TRACTABLE, VALOSSA, RESNAP, cortica

## AERIAL IMAGING AND MAPPING

Orbital Insight, planet, FEATUREX, CartoAware, VIDEO inform, terraloupe

## SOCIAL MEDIA, ENTERTAINMENT AND LIFESTYLE

OnCara, SECOND SPECTRUM, SOTERIA INTELLIGENCE, OLLY, Jukedeck, Neural Painting, VAULT, Sensifai, Sonalytic Systems, PIXONEYE

## AUDIO AND LANGUAGE AS A SERVICE

AISense, babblabs, CAPIO, CTI, gridshift, indico, KITT.AI, LEXALYTICS, MindMeld, MOD9, PAT, semanticmachines, volley, Captricity, LUMINOSO, AISPEECH, Hubino, intelliGo, MIND, Unisound, Bloomsbury AI, Cognicor, cortical.io, 1Checker, Intelligent Voice, LASTMILE, LEVERTON, aido, Replika, retechnica, spaCy, speech Matics, They Say, idio, Linguamatics, audioburst, BEYONDVERBAL

## HUMAN-MACHINE INTERFACE

neurala, neon, LIGHTHOUSE, SKINDROID, TERADEEP, gyrfalcon technology, ZERO AI, SOUND HOUND, IMOTIONS, mashgin, MOD9, KITT.AI, Algorithmic Intuition, keenresearch, trueface.ai, eyeris, SoundHound Inc., Sighthound, aNXI, iz.ai, AIBRAIN INC, babblabs, EMOTIBOT, MORPX, Mobvoi, AKA, Rokid, Face++ Cognitive Services, AISPEECH, VEO, intuition robotics, LIPSIGHT, emteq, eyeSight, IMAGRY, JUNGO, ORCAM, snips

## IT OPERATIONS AND SECURITY

CYLANCE, graphistry, sparkcognition, deep instinct, Cognitive ID, Cyberlytic, DARKTRACE

## DRONES AND ROBOTS

abundant ROBOTICS, CLEARPATH, Airware, ACCELERATED DYNAMICS, LILY, Iris Automation, UNIVERSAL, Qelzal, SKYDIO, OSARO, brain corp, AISPEECH, KINDRED, HOVER CAMERA, DLR, ACCELERATED DYNAMICS, SCORTEX, robart, evolve

## CRM AND HR

BICYCLE AI, Digital Genius, vufind, Eloquent Labs, MARAX AI, SentiSum, CYRA, AYLIEN, Celaton, re:infer

## SILICON PLATFORMS

REM, thinci, groq inc., Wave Computing, TENSTORRENT, gyrfalcon technology, brainchip, DEEP VISION, NOVUMIND, MYTHIC, Preferred Networks, Horizon Robotics, Cambricon, DEEPHi, Unisound, intelliGo, Think Force, Graphcore, AMOTIVE

## SECURITY, SURVEILLANCE, MONITORING

Deep Vision, DEEP SENTINEL, SHIELD AI, LIGHTHOUSE, NOVUMIND, gyrfalcon technology, cogniac, Sighthound, Igocian, PILOT.AI, athelas, SENTENAI, camio, trueface.ai, AiCure, SECOND SPECTRUM, algolux, senseme, Face++ Cognitive Services, DEEPHi, YITU, DEEPGLINT, intelliGo, intelLIfusion, Horizon Robotics, OpenCapacity, VISIOINGENII, KONUX, 3rdeye, Calipsa, alchera, xihelm, HOXTON, GETALERT, LIPSIGHT, FIFTH DIMENSION, anyvision, PointGrab, VIDEO inform

## AUTONOMOUS VEHICLES

QUANERGY, algolux, drive.ai, PlusAI, PITSTOP, AutoX, nauto, nuro, Velodyne LiDAR, ZOOX, ARGO AI, nuTonomy, DEEPSCALE, AURO, tu simple, AISPEECH, senseme, MINIEYE, LEAPMIND, netradyne, BLUE VISION, machines with vision, AMOTIVE, RoboCV, Mapillary, OXBOTICA, i4drive, Brodmann, FIVE AI, cognata, nexar, COGNITIVE PILOT, Morpheus Labs, SafeCue

## OPERATIONS, LOGISTICS AND MANUFACTURING

Citrine, clearmetal, Descartes Labs, DroneDeploy, SIGHT MACHINE, INVENIA, MAANA, Intellegens, micropsi industries, Terrabotics, AUGURY, PRESENSO, trax, 3DSignals

speak your mind