# About the Chamberlain Group

**CHAMBERLAIN GROUP**

LiftMaster    CHAMBERLAIN    myQ    tend

CPSG    SYSTEMS    Merlin    GRIFCO

Chamberlain Group (CGI) is a global leader in access solutions and products.

**Over 8,000 Employees Worldwide**

CGI is a global team with solutions and operations designed to serve customers in a variety of markets worldwide.

**VISION**

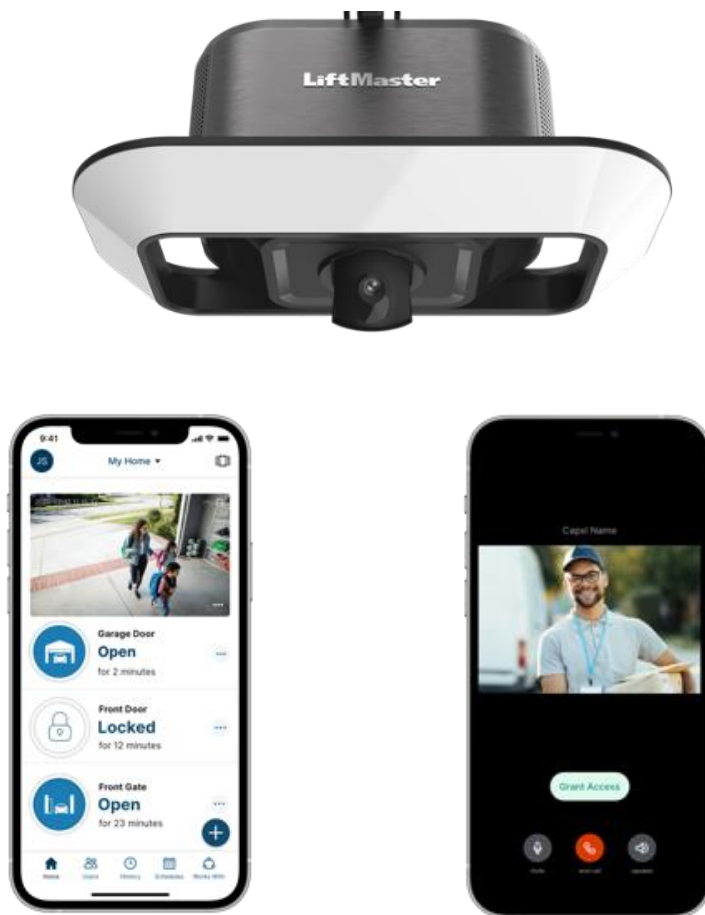Giving the Power of Access and Knowledge

**MISSION**

People everywhere rely on CGI to move safely through their world, confident that what they value most is secure within reach.

**END-MARKETS SERVED**

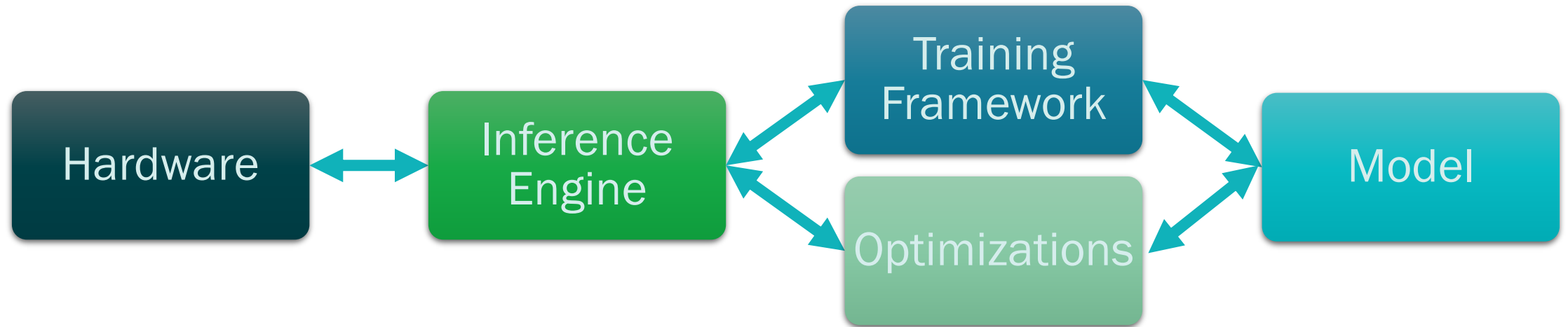Residential
Commercial
Automotive

**CHAMBERLAIN GROUP**

# About the Chamberlain Group

# More is possible than you think!

- Survey the landscape of edge inference implementation

- Explore software & hardware choices

- Examine model & optimization choices

# Everything is intertwined.



Hardware ↔ Inference Engine ↔ Training Framework / Optimizations ↔ Model

Hardware → Inference Engine ↔ Training Framework / Optimizations ↔ Model

Coral → TFLite → TensorFlow / 8-bit Quantized / Limited Layer Choice → Only Some Models

CHAMBERLAIN GROUP

# Hardware

# Hardware Choices

+GPU

**CPU**

**Cortex-A, x86**

# Hardware Choices

# Hardware Choices

| +SIMD | +DSP | +NPU | +FPGA | | +SIMD | +DSP | +GPU | +NPU | +FPGA |
|-------|------|------|-------|--|-------|------|------|------|-------|

| MCU Cortex-M, RISC-V | CPU Cortex-A, x86 | FPGA |
|----------------------|-------------------|------|

**CHAMBERLAIN GROUP**

# Hardware Choices

| +SIMD | +DSP | +NPU | +FPGA | +SIMD | +DSP | +GPU | +NPU | +FPGA |
|---|---|---|---|---|---|---|---|---|

| MCU Cortex-M, RISC-V | CPU Cortex-A, x86 | FPGA |
|---|---|---|

Multiple Cores & Clock Speed

On-Chip Memory, L1 & L2 Caches

Bus Speed

Register Count

Out-of-Order Execution

"Hidden" Fundamentals

# Hardware: Common Configurations

## Low Volume



USB



Software → Hardware

## Mass Production



MIPI



Memory,
IO, etc.

Hardware → Software

**CHAMBERLAIN GROUP**

# Model & Software

# Model Optimization Workflow

**Reference Model**

**Optimization**

**Convert or Compile**

**Deploy and Test**

- State-of-the-art
- FP32 GPU
- Purpose: validate your training data

- Choose backbone
- Choose head
- Quantization
- Pruning
- NAS

- Convert
- Compile

- Verify
- Review Choices
- Analyze
- Optimize
- Integrate

# Model Optimization Workflow

**Reference Model** → **Optimization** → **Convert or Compile** → **Deploy and Test**

**Reference Model**
- State-of-the-art
- FP32 GPU
- Purpose: validate your training data

**Optimization**
- Choose backbone
- Choose head
- Quantization
- Pruning
- NAS

**Convert or Compile**
- Convert
- Compile

**Deploy and Test**
- Verify
- Review Choices
- Analyze
- Optimize
- Integrate

CHAMBERLAIN GROUP

- Your problem space is different!

- Smaller datasets are usually OK
  - Smaller models need less data
  - Fine tuning needs less data

Research Datasets

Your Application

# Model Optimization Workflow

| Reference Model | Optimization | Convert or Compile | Deploy and Test |
|---|---|---|---|

**Reference Model**
- State-of-the-art
- FP32 GPU
- Purpose: validate your training data

**Optimization**
- Choose backbone
- Choose head
- Quantization
- Pruning
- NAS

**Convert or Compile**
- Convert
- Compile

**Deploy and Test**
- Verify
- Review Choices
- Analyze
- Optimize
- Integrate

**CHAMBERLAIN GROUP**

# Model Mash-Up

Output

Layer 5

Layer 4

Layer 3

Layer 2

Layer 1

Input

**Head & Neck:**
- Interprets results
- Inexpensive to fine-tune
- Lower data requirements than backbone

**Backbone:**
- Extracts features
- Costly to train; needs lots of data & time
- **Recommendation: pre-trained weights**

# Optimization Options

Novice ➡ Amateur ➡ Professional ➡ Expert

| | Novice | Amateur | Professional | Expert |
|---|---|---|---|---|
| **Model Structure** | Model Zoo | Community-Supported | Mix & Match Head & BB / Code it yourself | |
| **Model Optimizations** | Pre-Optimized Model | Quantization, Pruning, Compression | | Decomposition / NAS / OFA |
| **Inference Engine** | Pre-Optimized Runtime | Community-Supported Runtimes | | Hand-Optimized Code |
| **Training Data** | Pretrained | Fine-tune with your data | | Train from scratch |

Deeplite Commercial ← → Open Source
EDGE IMPULSE

**CHAMBERLAIN GROUP**

# Model Optimization Workflow

**Reference Model** → **Optimization** → **Convert or Compile** → **Deploy and Test**

**Reference Model**
- State-of-the-art
- FP32 GPU
- Purpose: validate your training data

**Optimization**
- Choose backbone
- Choose head
- Quantization
- Pruning
- NAS

**Convert or Compile**
- Convert
- Compile

**Deploy and Test**
- Verify
- Review Choices
- Analyze
- Optimize
- Integrate

**CHAMBERLAIN GROUP**

- **Runtime**
  - Model is interpreted
  - Model deployed separately
  - Easier OTA updates

- **Compiler**
  - Model is compiled
  - Model is part of firmware
  - Weights are often constants

- **Code Optimizations**
  - Memory Usage
    - Cache-aware (e.g., tiling)
    - Efficient register usage
  - Vectorization
    - Use SIMD
    - Use DSP, NPU, GPU
  - Parallelization

- **Consensus over time**
  - No model gets it right all the time
- **High frame rate:**
  - More samples for consensus
  - Lower per-sample accuracy
- **Low frame rate:**
  - Fewer samples for consensus
  - Higher per-sample accuracy

**100 ms inference time does NOT mean 10 FPS!**

**Reserve CPU cycles for:**

- Ingesting from the sensor/buffer
- Interpreting the output
- Network
- Other app functions
- Temperature management

# Example

# Example: Object Detection on ArmV7

| Task | Vehicle Detection |
|------|-------------------|
| **Reference Model** | YOLOv5-s, FP32, PyTorch |
| **Compute Constraints** | ARMv7 w/NEON, no accelerators |

Seeed NPI i.MX6ULL
Cortex A7 @ 800 MHz

Raspberry Pi 2 B v1.1
Cortex A7 @ 900 MHz

ASUS Tinkerboard (v1)
Cortex A17 @ 1.8 GHz

CHAMBERLAIN GROUP

# Optimization Options

## Single-Core Inference Time (Millisecond)



| Model | Seeed NPI i.MX6ULL Cortex A7 800 MHz | Raspberry Pi 2 B v1.1 Cortex A7 900 MHz | ASUS TinkerBoard v1 Cortex A17 1.8 GHz |
|---|---|---|---|
| Pelee (NCNN) | 5476 | 4605 | 1163 |
| MobileNetV2-YOLO (NCNN) | 3500 | 3210 | 691 |
| NanoDet (NCNN) | 680 | 609 | 104 |
| QuickYOLO (LCE) | 637 | 542 | 111 |
| Xailient | 67 | 59 | 14 |

**Legend:**
- Seeed NPI i.MX6ULL Cortex A7 800 MHz
- Raspberry Pi 2 B v1.1 Cortex A7 900 MHz
- ASUS TinkerBoard v1 Cortex A17 1.8 GHz

# Models & Tools Tested

## Single-Core Inference Time (Millisecond)



**Pelee (NCNN)**
- 5476
- 4605
- 1163

**MobileNetV2-YOLO (NCNN)**
- 3500
- 3210
- 691

**NanoDet (NCNN)**
- 680
- 609
- 104

**QuickYOLO (LCE)**
- 637
- 542
- 111

**Xailient**
- 67
- 59
- 14

0   1000   2000   3000   4000   5000   6000

- ■ Seeed NPI i.MX6ULL Cortex A7 800 MHz
- ■ Raspberry Pi 2 B v1.1 Cortex A7 900 MHz
- ■ ASUS TinkerBoard v1 Cortex A17 1.8 GHz

### Why so much faster?
- Both 32-bit ARMv7
- NEON: 64-bit vs 32-bit
- FP: 16 registers vs 32 registers
- Out-of-order execution
- Deeper pipeline
- DMIPS/MHz: 4.0 vs 1.9

**CHAMBERLAIN GROUP**

# Models & Tools Tested

## Single-Core Inference Time (Millisecond)



| Model | Seeed NPI i.MX6ULL | Raspberry Pi 2 B v1.1 | ASUS TinkerBoard v1 |
|---|---|---|---|
| Pelee (NCNN) | 5476 | 4605 | 1163 |
| MobileNetV2-YOLO (NCNN) | 3500 | 3210 | 691 |
| NanoDet (NCNN) | 680 | 609 | 104 |
| QuickYOLO (LCE) | 637 | 542 | 111 |
| Xailient | 67 | 59 | 14 |

Legend:
- Seeed NPI i.MX6ULL Cortex A7 800 MHz
- Raspberry Pi 2 B v1.1 Cortex A7 900 MHz
- ASUS TinkerBoard v1 Cortex A17 1.8 GHz

**CHAMBERLAIN GROUP**

## Single-Core Inference Time (Millisecond)



| | | Seeed NPI i.MX6ULL (Cortex A7 800 MHz) | Raspberry Pi 2 B v1.1 (Cortex A7 900 MHz) | ASUS TinkerBoard v1 (Cortex A17 1.8 GHz) | Speed improvement |
|---|---|---|---|---|---|
| | Pelee (NCNN) | 5476 | 4605 | 1163 | 6x |
| | MobileNetV2-YOLO (NCNN) | 3500 | 3210 | 691 | 10x |
| 80-class | NanoDet (NCNN) | 680 | 609 | 104 | 50x |
| 80-class | QuickYOLO (LCE) | 637 | 542 | 111 | 54x |
| 1-class | Xailient | 67 | 59 | 14 | 509x |

Speed improvement compared to YOLOv5-s on PyTorch

■ Seeed NPI i.MX6ULL Cortex A7 800 MHz
■ Raspberry Pi 2 B v1.1 Cortex A7 900 MHz
■ ASUS TinkerBoard v1 Cortex A17 1.8 GHz

## Model Accuracy (Custom Vehicle Dataset)



| Model | AP50 | mIoU | Speed Improvement |
|---|---|---|---|
| YOLOv5-S (PyTorch) — Reference Model | 99% | 95% | |
| NanoDet (NCNN) | 95% | 94% | 50x |
| QuickYOLO (LCE) | 90% | 83% | 54x |
| Xailient | 90% | 79% | 509x |

Legend: ■ AP50  ■ mIoU

**CHAMBERLAIN GROUP**

# Conclusions

- Inference on edge devices has become both possible and practical

- Small hardware features can make a big difference in speed

- Selecting the right model and the right inference engine for your hardware can expand the scope of what is possible

**Detectors used in the Example:**

YOLOv5

https://github.com/ultralytics/yolov5

NanoDet

https://github.com/RangiLyu/nanodet

QuickYOLO

https://github.com/tehtea/QuickYOLO

Xailient

https://www.xailient.com/

**Chamberlain Group:**

https://chamberlaingroup.com

**Note:**

Many more links and resources are available at the end of the slide deck.

# Backup Material

# Half of implementing deep learning is fighting Python & C++ errors and resolving library incompatibilities.

Pay close attention to documented versions!

Use "virtualenv"

Become a CMAKE expert!

# Papers

| Paper Title | URL |
|---|---|
| Larq Compute Engine: Design, Benchmark, and Deploy State-of-the-Art Binarized Neural Networks | https://arxiv.org/abs/2011.09398 |
| Latent Weights Do Not Exist: Rethinking Binarized Neural Network Optimization | https://arxiv.org/abs/1906.02107 |
| FCOS: Fully Convolutional One-Stage Object Detection | https://arxiv.org/abs/1904.01355 |
| Bridging the Gap Between Anchor-based and Anchor-free Detection via Adaptive Training Sample Selection | https://arxiv.org/abs/1912.02424 |
| Generalized Focal Loss: Learning Qualified and Distributed Bounding Boxes for Dense Object Detection | https://arxiv.org/abs/2006.04388 |
| ShuffleNet V2: Practical Guidelines for Efficient CNN Architecture Design | https://arxiv.org/abs/1807.11164 |

# Edge Inference Engines
## Generic, Open-Source

| Inference Engine | Type | Notes | |
|---|---|---|---|
| TFLite | Runtime | Runtime for TF/Keras | https://www.tensorflow.org/lite |
| TFLite Micro | Runtime | TFLite for MCUs | https://www.tensorflow.org/lite/microcontrollers |
| Larq Compute Engine | Runtime | Binarized TFLite | https://github.com/larq/compute-engine |
| NCNN | Runtime | Tencent runtime | https://github.com/Tencent/ncnn |
| MNN | Runtime | Alibaba runtime | https://github.com/alibaba/MNN |
| Apache TVM | Compiler | Compiler and optimizer | https://tvm.apache.org/ |
| Apache MicroTVM | Compiler | TVM for MCUs | https://tvm.apache.org/docs/microtvm/index.html |
| Glow | Compiler | Compiler for ONNX | https://ai.facebook.com/tools/glow/ |
| Microsoft ELL | Compiler | Compiler | https://github.com/Microsoft/ELL |
| deepC | Compiler | ONNX -> LLVM | https://github.com/ai-techsystems/deepC |
| NNoM | Library | Keras -> C | https://github.com/majianjia/nnom |

*Note: This list is not comprehensive.*

**CHAMBERLAIN GROUP**

# Edge Inference Engines
## Hardware-Specific and Commercial

| Vendor | Inference Engine | Type | Notes |
|--------|------------------|------|-------|
| Nvidia | TensorRT | Runtime | Support for Nvidia GPUs, such as Jetson Nano |
| Arm | Arm® NN | Runtime | Optimized for Arm Cortex-A CPU, Mali GPU, Ethos NPU |
| Arm | CMSIS-NN | Library | Library used by various runtimes and compilers |
| NXP | NXP eIQ™ | Both | Optimized for NXP; Tflite and Glow with Arm-NN & CMSIS-NN |
| Qualcomm | SNPE | Runtime | For Qualcomm Snapdragon processors |
| Intel | OpenVINO™ | Runtime | Runtime for Intel products, including Movidius |
| Morpho | SoftNeuro | Runtime | Commercial platform; limited details publicly available. |
| Edge Impulse | EON | Compiler | Commercial platform; Targeted at Microcontrollers |
| STMicro | STM32Cube.AI | Compiler | Optimized for STM32 |
| Kendryte | nncase | Compiler | Kendryte K210; https://github.com/kendryte/nncase |

*Note: This list is not comprehensive.*

# Model Optimization Tools
## Generic, Open-Source

| Tool | Framework(s) | URL |
|------|-------------|-----|
| TensorFlow MOT | TensorFlow | https://www.tensorflow.org/model_optimization |
| Microsoft NNI | PyTorch | https://github.com/microsoft/nni |
| IntelLabs Distiller | PyTorch | https://github.com/IntelLabs/distiller |
| Riptide | TensorFlow + TVM | https://github.com/jwfromm/Riptide |
| Qualcomm AIMET | PyTorch, TensorFlow | https://github.com/quic/aimet |

*Note: This list is not comprehensive.*

**CHAMBERLAIN GROUP**

# Model Optimization Tools
## Hardware-Specific and Commercial

| Tool | Framework(s) | URL |
| --- | --- | --- |
| OpenVINO NNCF | PyTorch | https://github.com/openvinotoolkit/nncf |
| NXP eIQ | TensorFlow, TFLite, ONNX | https://www.nxp.com/design/software/development-software/eiq-ml-development-environment:EIQ |
| Deeplite | PyTorch, TensorFlow, ONNX | https://www.deeplite.ai/ |
| Edge Impulse | Keras | |

*Note: This list is not comprehensive.*

CHAMBERLAIN GROUP

# Peripheral Accelerators

| Product | Off-the-Thelf SBC | USB |
|---|---|---|
| Nvidia GPU | Jetson Nano, TX1, TX2 | - |
| Movidius Myriad X | - | Intel Neural Compute Stick 2 |
| Google Edge TPU | Coral Dev Board, Dev Board Mini | Coral USB Accelerator |
| Gryfalcon Lightspeeur® | - | Orange Pi AI Stick Lite |
| Rockchip RK1808 | - | Toybrick RK1808 |

*Note: This list is not comprehensive.*

# SoCs w/Embedded Accelerators

| Product | Acceleration | Single Board Computer |
|---------|-------------|----------------------|
| Qualcomm Snapdragon (various) | DSP + GPU (+NPU) | *(by request only)* |
| Ambarella CV2, CV5, CV22S, CV25S, CV28M | DSP + NPU | *(by request only)* |
| NXP i.MX 8 | DSP + GPU | SolidRun $160+ |
| NXP i.MX 8M Plus | DSP + GPU + NPU | SolidRun, Wandboard $180+ |
| Rockchip RK3399Pro | NPU | Rock Pi N10 $99+ |
| Allwinner V831 | NPU | Sipeed MAIX-II Dock $29 |
| Sophon BM1880 | NPU | Sophon Edge $129 |

*Note: This list is not comprehensive.*

# MCUs for Inference

| Vendor | Product | Features that support inference |
|---|---|---|
| Various | Cortex-M4/7/33/35P | SIMD instructions, FPU; Future Ethos-U55 microNPU |
| Raspberry Pi | RP2040 | Memory, bus fabric |
| Maxim Integrated | MAX78000 | Cortex-M4, CNN accelerator |
| Kendryte | K210 | DNN accelerator |
| Espressif | ESP32-S3 | SIMD instructions, FPU |

*Note: This list is not comprehensive.*

**CHAMBERLAIN GROUP**