

Embedding Programmable DNNs in Low Power SOCs



Steve Teig 22nd May 2018

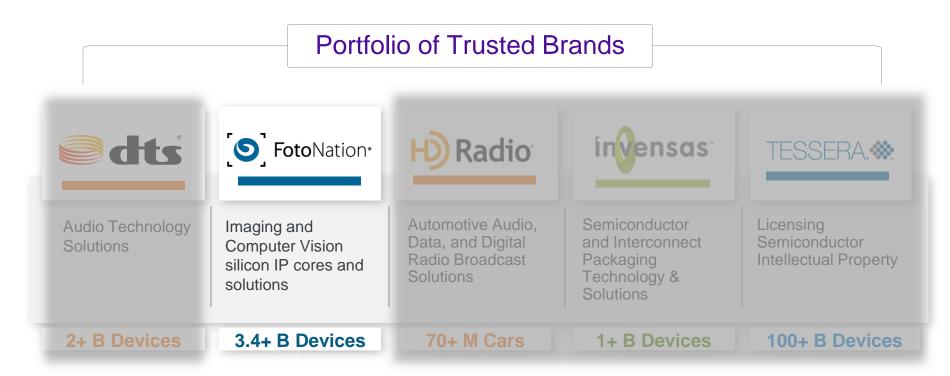


Company Overview



FotoNation – trusted brand of XPERI







What, How and Why



What – Al Enabled Imaging at the Edge





Always-on inference: operates even while the device is "off"

(e.g., ultra low power FD/FR as an enabler)



Driver State
Sensing for
autonomous driving

(e.g., always ON driver assistant)



Smart IoT: TVs to drones to microwave ovens to ...

(e.g., ultra low power people detection)

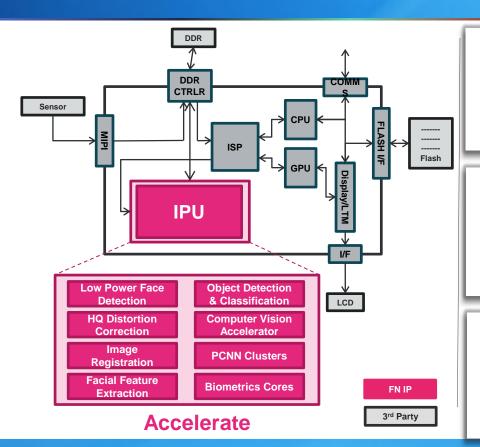


Head-mounted displays for AR or MR

(e.g., ultra low power IRIS, eye gaze, scene understanding,, etc..)

How - Image Processing Unit hybrid deployment





Understand



- Face, People, Object Detection, Segmentation & Tracking
 • Scene Classification and
- recognition

Enhance



- De-warping & Stitching Stabilization with Rolling Shutter Correction
- HDR & LTM

Personalize



- Visible and NIR 3D FR
- IRIS recognition Liveliness Detection & Continuous recognition (hand jitter, facial, etc.)



Why - Advantages



- Ultra high performance for always ON features via dedicated accelerators (e.g., facial detection and recognition, resampling, LTM, etc.)
- High flexibility and configurability after silicon via programmable inference engines
- Flexibility and configurability for various markets via the concept of clusters
- Unmatched performance and power consumption for concurrent understanding, enhancing & personalization imaging features

Examples



UNDERSTAND



PERSONALIZE





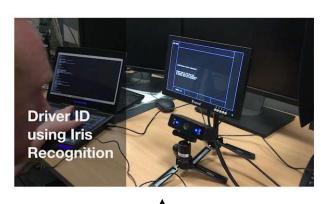


IMAGE PROCESSING UNIT (IPU)



What is the IPU?



IPU in a Nutshell



Collection of configurable RTL engines connectable to preprocess and analyze incoming image and video streams

PCNNs

- PCNNs for reconfigurable functionality
- Concurrent support of multiple real-time networks

Dedicated Cores

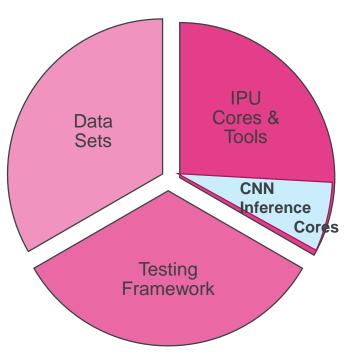
- Facial & people analytics (AI/ML)
- Stabilization (HQ resampling, analytics)
- Optional Depth (AI/ML)

Pre-processing

- Multi-resolution stream generation
- Local tone mapping enhancement (significantly improves detection ratios)
- High-quality, low-latency distortion correction

IPU – the supporting ecosystem



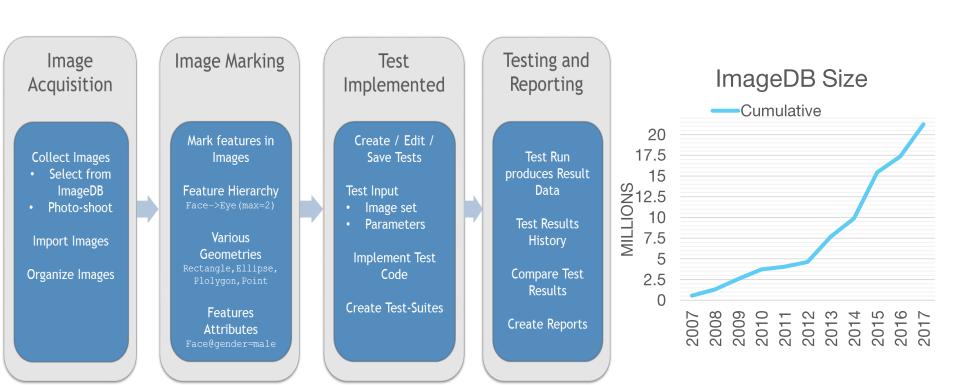


... inference cores are a small part of the total

- IPU
 - RTL (image pre-processing, dedicated analytics and inference cores)
 - Tools for programming, training and debugging
- Testing Framework (a.k.a. ImageDB) test framework supporting acquisition, marking, testing and reporting
- Data Sets (a.k.a. CV Infra) computer vision infrastructure supporting 2D and 3D image sets acquisition, annotation, marking and training sets generation with ground truth

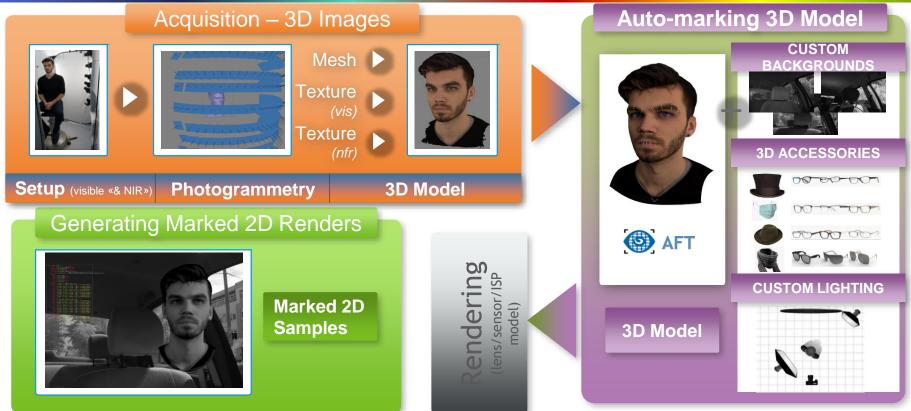
Testing Framework (a.k.a. Image DB)





Computer Vision Infrastructure







PCNNs and PCNN-Clusters



FotoNation's Programmable CNN Engine (PCNN)



Built in preprocessing on the fly imaging engine for (e.g., layer 0)

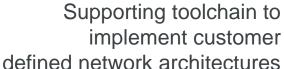
Local memory for fast data access and reduced memory BW Support for compression, quantization, decryption



PCNN 1.2 - 36 MAC/cycle or

PCNN 2.0 - 512 MAC/cycle

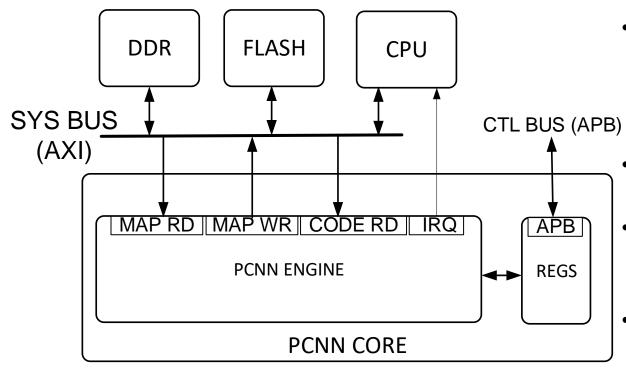
Designed for very low latency and real-time network inference





FotoNation's PCNN Highlights



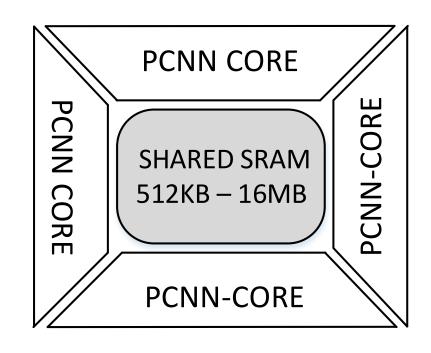


- Low power consumption (22n FDSOI tech)
 - 18 mW for PCNN 1.2
 - 120 mW for PCNN 2.0
- Built in real-time, on-the-fly NN decryption engine
- Separate memory channels for network fetch and intermediate layers/input
- Separate cache for code/net & intermediate layers

PCNN-Cluster Concept



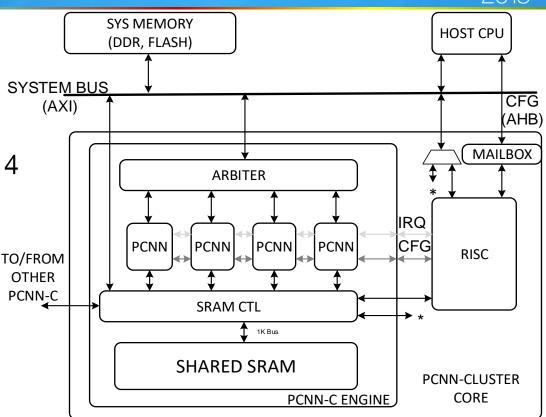
- PCNN–C accelerates any type of CNN
- Highly configurable and scalable architecture
- Multiple PCNN IP cores connected to a shared common memory (e.g., SRAM)
- Configurable with up to 4 PCNN cores



PCNN-Cluster Highlights



- Clusters can have a mix of PCNN 2.0 or 1.2 cores
- Can execute single or up to 4 individual networks at the same time
- Several clusters can be chained for increased computing power and capabilities

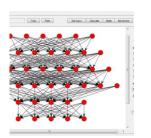


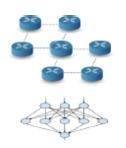
Security and Protection – the Danger













Networks sit in the device's permanent storage

Device can be accessed and storage contents read

Neural processor makers offer network transfer tools

Networks representation patterns can be identified and localized in the storage contents

Network architecture & weights extraction

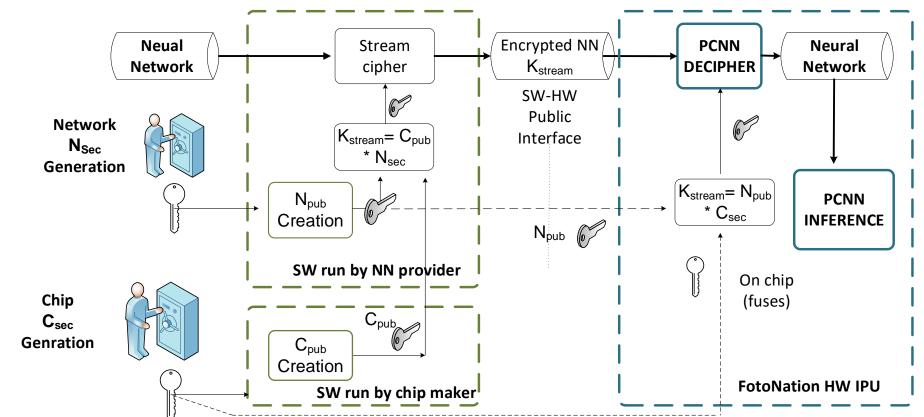
Once the network representation is known, architecture and weights values can be obtained

Network re-map and inference on alternative architectures

Networks can be remapped and inferred on alternative architectures as own

Solution







DBI® 3D AI Enabled Chip Architectures using IPU

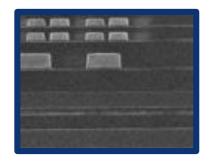


Digital Bond Interconnect platforms for 3D semiconductors – DBI®



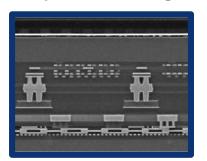
ZiBond®

Homogeneous Bonding



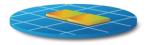
Scalable to 1,000,000 interconnects per mm² DBI®

Hybrid Bonding





Wafer to Wafer (W2W)



Die to Wafer (D2W)



Die to Die (D2D)

Technology
Development &
Optimization



3D Design & Architecture



Simulation



Materials Characterization



Wafer/Die Bonding & Processing



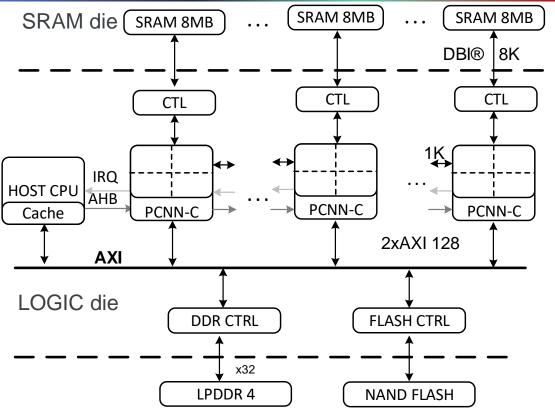
Reliability



Failure Analysis

Al Enabled Chip Architecture with IPU & DBI®







- Array of 2µm pitch DBI®
 interconnects (between silicon dies) up to 250,000 vertical interconnects
 per mm² enables groundbreaking
 computing architecture
- Very short vertical interconnects offer ultra high performance at very low power



Conclusions



Conclusions



- Edge computing requires <u>low power and real-time</u> response
- Huge investment <u>beyond inference cores</u> in development, data preparation, testing and validation of the algorithms
- Edge computing brings <u>IP security complications</u> with regard to reverse engineering
- Edge computing is recommended to have real time, on-the fly network <u>protection</u> with no meaningful compromise to latency, area, and power consumption; FotoNation Cores do
- Edge intelligence with IPU is the way forward: for <u>power efficiency</u>, <u>latency and privacy/security reasons</u>

XPERI enables ultimate user experience making possible AI enabled imaging in the device – a.k.a. the Edge ...

Resources



- DBI®: https://www.invensas.com/solutions/waferdie-bonding/
- IPU: https://www.fotonation.com/products/image-processing-unit/
- Deep Learning course: https://www.coursera.org/specializations/deep-learning#courses
- Deep Learning book: http://www.deeplearningbook.org/



Thank you!





Additional Slides



PCNN-C performance



PCNN-C Configuration scenarios

	PCNN-C Lite	PCNN-C Base	PCNN-C Performance	PCNN-C Future
Number of PCNN 2.0 Cores	0	1	2	4
Number of PCNN 1.2 Cores	2	3	2	0

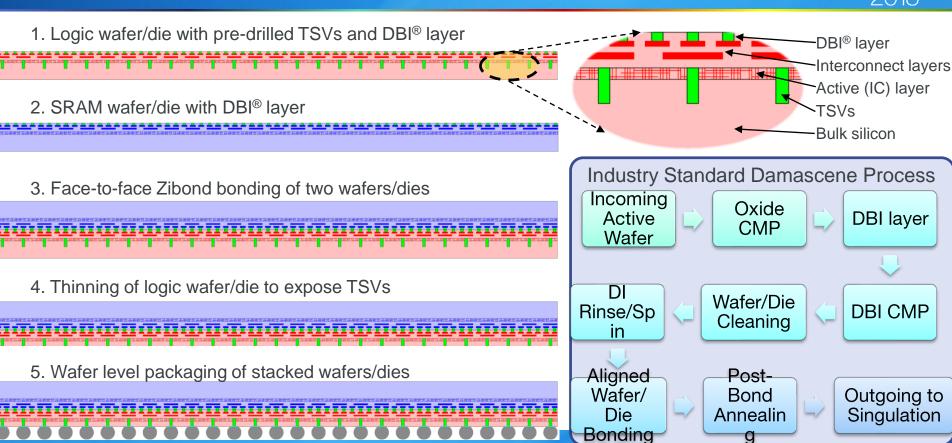
PCNN-C Performance examples on various networks

	PCNN-C Lite	PCNN-C Base	
Yolov2 [FPS]	30	250	
Body Pose [FPS]	ı	15	
CCOT [FPS]	450	4700	
SEGNET [FPS]	10	80	
AlexNet [FPS]	60	400	
LeNet [FPS]	20000	100000	
GoogleNetVx [FPS]	30	270	
VGG16 [FPS]	3	12	
FotoNation BD [FPS]	44	300	
FotoNation FR (largest			
network on visible)	44	300	
[FPS]			



DBI® explained





DBI® with Zibond: the most advanced 3D interconnects



