



HAL
open science

NimbleAI: Towards Neuromorphic Sensing-Processing 3D-integrated Chips

Xabier Iturbe, Nassim Abderrahmane, Jaume Abella, Sergi Alcaide, Eric Beyne, Henri-Pierre Charles, Christelle Charpin-Nicolle, Lars Chittka, Angélica Dávila, Arne Erdmann, et al.

► **To cite this version:**

Xabier Iturbe, Nassim Abderrahmane, Jaume Abella, Sergi Alcaide, Eric Beyne, et al.. NimbleAI: Towards Neuromorphic Sensing-Processing 3D-integrated Chips. DATE '23 Design, Automation and Test in Europe Conference, Apr 2023, Antwerp, Belgium. pp.6, 10.23919/DATE56975.2023.10136952 . hal-04412850

HAL Id: hal-04412850

<https://hal.science/hal-04412850>

Submitted on 4 Mar 2024

HAL is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

NimbleAI: Towards Neuromorphic Sensing-Processing 3D-integrated Chips

Xabier Iturbe, Nassim Abderrahmane, Jaume Abella, Sergi Alcaide, Eric Beyne, Henri-Pierre Charles, Christelle Charpin, Lars Chittka, Angélica Dávila, Manil Dev Gomony, Arne Erdmann, Carles Estrada, Ander Fernández, Anna Fontanelli, José Flich, Alejandro Hernán Gloriani, Radu Grosu, Carles Hernández, Daniele Ielmini, Eric Isusquiza, David Jackson, Maha Kooli, Nicola Lepri, Bernabé Linares-Barranco, Jean-Loup Lachese, Martxel Lasa, Eric Laurent, Menno Lindwer, Frank Linsenmaier, Mikel Luján, Karel Masařík, Nele Mentens, Orlando Moreira, Luca Peres, Jean-Philippe Noel, Arash Pourtaherian, Christoph Posch, Peter Priller, Zdenek Prikryl, Felix Resch, Oliver Rhodes, Todor Stefanov, Moritz Storing, Sander Stuijk, Michele Taliercio, Marcel van de Burgwal, Geert van der Plas, Elisa Vianello, and Pavel Zaykov

Abstract—This article introduces the Horizon Europe *NimbleAI* research project, which brings together 19 EU and UK partners covering most of the semiconductor value chain. *NimbleAI* leverages key principles of energy-efficient visual sensing and processing in biological eyes and brains, and harnesses the latest advances in 3D stacked silicon integration, to create an integral sensing-processing neuromorphic architecture that efficiently and accurately runs computer vision algorithms in area-constrained endpoint chips. The rationale behind the *NimbleAI* architecture is: sense data only with high information value and discard data as soon as they are found not to be useful for the application (in a given context). The *NimbleAI* sensing-processing architecture is to be specialized after-deployment by tuning system-level trade-offs for each particular computer vision algorithm and deployment environment. The objectives of *NimbleAI* are: (1) 100x performance per mW gains compared to state-of-the-practice solutions (i.e., CPU/GPUs processing frame-based video); (2) 50x processing latency reduction compared to CPU/GPUs; (3) energy consumption in the order of tens of mWs; and (4) silicon area of approx. 50 mm².

Index Terms—Neuromorphic, computer vision, 3D silicon, event-based vision, in-memory computing, eFPGA, RISC-V, virtual neural networks, light-field vision, online learning

X. Iturbe, C. Estrada, A. Fernández, M. Lasa and A. Dávila are with IKERLAN, Basque Country (Spain); C. Hernandez and J. Flich are with Universitat Politècnica de València, Spain; J. Abella and S. Alcaide are with Barcelona Supercomputing Center (BSC), Catalonia (Spain); F. Resch and R. Grosu are with TU Wien, Austria; G. Van der Plas, M. van de Burgwal, M. Storing and E. Beyne are with IMEC, Belgium; A. Erdmann is with Raytrix, Germany; E. Isusquiza is with ULMA Medical Imaging, Basque Country (Spain); N. Abderrahmane, J.L. Lachese and E. Laurent are with MENTA, France; N. Mentens and T. Stefanov are with Universiteit Leiden, Netherlands; P. Zaykov, Z. Prikryl and K. Masařík are with CODASIP, Czech Republic; M. Lindwer, O. Moreira and A. Pourtaherian are with GRAI Matter Labs (GML), Netherlands; N. Lepri and D. Ielmini are with Politecnico Milano, Italy; O. Rhodes, M. Luján, L. Peres and D. Jackson are with University of Manchester, UK; B. Linares-Barranco is with CSIC, Spain; A. Fontanelli and M. Taliercio are with Monozukuri (MZ Technologies), Italy; M. Kooli, H.P. Charles, C. Charpin and J.P. Noel are with CEA-LIST, France; E. Vianello is with CEA-LETI, France; P. Priller is with AVL List, Austria; L. Chittka is with Queen Mary University of London, UK; A.H. Gloriani and F. Linsenmaier are with Viewpointssystem, Austria; C. Posch is with PROPHESEE, France; M. Gomony and S. Stuijk are with TU Eindhoven, Netherlands. All co-authors are listed in alphabetical order except for the main and corresponding author (xiturbe@ikerlan.es).

I. INTRODUCTION

Processor architectures (e.g., CPU/GPUs) are very inefficient compared to the biological eye-brain system, which is honed by natural selection and applies the fundamental energy-saving principle of capturing, processing and storing data only when necessary. Hence, eyes continuously sense and encode the changing surrounding environment in a way that is manageable for the brain.

The recently started *NimbleAI* project leverages key principles of energy-efficient light detection in eyes and visual information processing in brains to create an integral sensing-processing neuromorphic chip that adopts the biological data economy principle at different system levels, and builds upon the latest advances in 3D stacked silicon integration. *NimbleAI* aims to deliver two world’s firsts: (1) a light-field dynamic vision sensor for monocular image-based depth perception; and (2) an event-driven end-to-end perception stack (*‘visual pathway’*) that runs industry standard Convolutional Neural Networks (CNNs). Since manufacturing a full 3D testchip is prohibitively expensive, *NimbleAI* will prototype key components via small-scale 2D stand-alone testchips. This cost-effective use of silicon allows us to produce high confidence research conclusions and silicon-proven neuromorphic IP.

This article discusses the main functioning principles of the *NimbleAI* architecture and existing system-level trade-offs: (1) sense only significant light changes (visual events) at the optimal spatio-temporal resolution; (2) distill sensed visual events to increase information-efficiency; (3) process selected information-rich events using minimal energy at the optimal DVFS point; and (4) route event-flows across the 3D stacked sensing-processing architecture to minimize data movement along shortest physical paths. The article is organized as follows. Section II introduces the major challenges of AI-enabling technologies that are addressed by *NimbleAI*. Section III outlines the overall *NimbleAI* concept and section IV describes the proposed architecture. Finally, Section V sums up the main takeaways to conclude the paper.

II. CURRENT CONTEXT AND CHALLENGES

NimbleAI deals with four main interrelated challenges and limitations of current AI-enabling technologies.

C1.- Complexity of AI models: Accuracy of computer vision algorithms is commonly opposed to efficiency [1]. CNNs are typically scaled up to increase accuracy by adding more layers or by enlarging these to process images at a higher resolution. On the other hand, state-of-the-practice edge CNNs typically rely on downscaling the resolution of full input images to keep workloads manageable by current inefficient processing architectures (see C3), thus sacrificing accuracy. Inaccuracies become greater when shrinking large industry standard CNNs to fit in resource-constrained edge and endpoint devices.

C2.- Performance and latency: State-of-the-practice computer vision systems are frame-based, which means that they periodically acquire and process full-size images in a layer-after-layer mode. Hence, the computation of one layer must be completed on the whole frame before the computation of the next layer starts. This results in growing inference delays as algorithms include more layers and sensor resolution increases.

C3.- Energy-efficiency of processor architectures: The current state-of-the-practice processor landscape includes general-purpose (CPU/GPU) and AI-specialized (NPU/TPU) architectures. CPU/GPUs are largely inefficient due to the continuous back-and-forth transfers of data (and instructions) with memory [2], whereas efficiency improvements brought about by NPUs (Neural Processing Units) and TPUs (Tensor Processing Units) depend to a great extent on the ability of the host CPU to split AI processing into matrix operations of similar dimensions to those for which the NPU/TPU architecture was optimized [3]. State-of-the-art neuromorphic architectures, on the other hand, implement brain-inspired (event-driven) neural networks to enormously increase energy-efficiency as they process only changes in their inputs [4]. Yet, only a few neuromorphic architectures promise to meet the high energy-efficiency levels with low energy budgets required at the endpoint (e.g., Innatera, SynSense, GrAI Matter Labs - GML, etc.). An important limitation of neuromorphic chips is that the size of neural networks that can run is restricted by the implemented neuron count in silicon. Innatera and Synsense chips implement only 1,000 neurons, greatly limiting their use to one dimensional applications such as audio. On the other hand, TrueNorth is the largest chip that IBM has ever built: at 500 mm² can hold only 1 M neurons [5] while real-world (image) applications typically require 10-20 M neurons and endpoint chips are typically 50 mm².

C4.- System integration: CPU/GPUs and NPU/TPUs are not typically integrated such that they can seamlessly and efficiently process data streams from sensors or interface to pre- and post-processing kernels. For example, TPUs do not have image sensor interfaces and hence need to rely on a host processor to capture and transmit video sequences to the TPU engine. For each video frame, this process may take factors more time than the TPU's actual AI processing of that same frame. Similar constraints hold for GPU and NPUs.

III. THE NIMBLEAI CONCEPT

NimbleAI considers that processing begins in the sensor. In fact, important efficiency gains are expected from the use of in-sensor analogue logic and novel Dynamic Vision Sensing (DVS) concepts that will be investigated for the first time in this project. These concepts include: (1) *digital-foveation* to dynamically allocate sensing resolution based on the information value brought about by each sensor region; and (2) coupling of light-field microlenses with the dynamic vision sensor to enable *event-based light-field perception*.

NimbleAI will study techniques to capture and optimally represent the spatio-temporal evolution of 3D scenes using minimal visual event-flows that match the optimization features implemented in the downstream processing and inference engines, and thus reduce energy consumption and latency of the whole architecture. We note that increasing the amount of meaningful information that can be obtained by scaling up DVS resolution is in fact one of the major open challenges in event-based vision [6]. The expectation is that by investing some computing power and energy to gain some situational awareness early, a major reduction of the amount of data to be processed will be achieved, saving lots of energy by doing that. This is inspired by unconscious visual processing and neural signalling in biological systems, and hence will be largely invisible to the user application yet adjustable through user-driven directives. It is also remarkable that NimbleAI aims to demonstrate an event-driven end-to-end pipeline that can run industry standard (large) AI models such as CNNs, to improve state-of-the-art approaches that are very limited (e.g., [7]).

As shown in Fig. 1, one of the novel system-level bio-inspired concepts that will be explored are event-driven *visual pathways* for optimal sensing and processing of feature-rich regions of interest (ROIs). We pose that visual pathways are an elegant way to answer challenge C4 and harness the increased bandwidth brought about by 3D integration, taking advantage of the irregular distribution of visual information and uneven temporal dynamics in the scene. Visual pathways will be assigned to ROIs in a one-to-one fashion: each pathway will span the sensor area determined by the assigned ROI and will include dedicated Through Silicon Viases (TSVs) to downstream sensed visual events to the processing and inference engines in the interior layers of the 3D stacked architecture. In addition to increasing bandwidth, 3D integration opens new opportunities to improve efficiency as the DVS resolution increases. Namely, events can be locally identified across the 3D structure reducing the number of ID bits needed. Each visual pathway will be configured independently and dynamically, from sensor to processing, at the accuracy (e.g., sensor resolution) and latency (e.g., DVFS settings) levels determined for that image region based on its dynamics and information value.

NimbleAI envisions a two-stage inference approach, where the two stages will reinforce each other to perform more efficiently as the deployment environments become more familiar and visual stimuli are better understood.

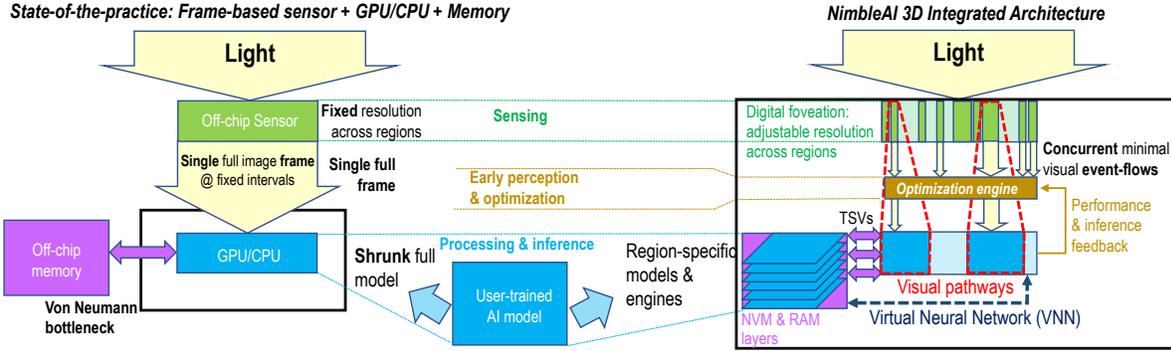


Fig. 1: NimbleAI approach (right) vs State-of-the-practice CPU/GPU (left).

1.- An always-on *early perception and optimization stage* implements selective attention algorithms to identify ROIs and configure accordingly the visual pathways. This includes selecting the most appropriate sensor resolution for each ROI and routing sensed visual events to the most appropriate user models (e.g., CNNs) running on downstream engines for efficient end-to-end region inference. NimbleAI will explore ultra-low energy and low-latency advantages of Spiking Neural Networks (SNNs) to power the early perception and optimization stage, as well as energy-efficient online learning rules to achieve specialization in dealing with deployment-specific visual stimuli. Hence, optimization SNNs will receive inference feedback from user models to continuously improve on dynamically selecting ROIs. This can be seen as a partial knowledge transfer from user-trained models to the early perception and optimization stage to optimize overall functioning. Hence, energy consumed to complete end-to-end inference also serves the purpose of adjusting the energy-saving mechanism in the early optimization stage. SNNs will also receive performance feedback and voltage/temperature analytics from monitors embedded along the visual pathways to continuously learn how to tune the execution conditions to be more efficient.

2.- An *inference stage* implements pre- and post-processing kernels on programmable hardware (i.e., MENTA eFPGA) and runs industry standard AI models on an event-driven dataflow accelerator (i.e., GML NeuronFlow) [8]. Processing in these components will be on-demand and optimized for the specific characteristics of each ROI. To deal with challenge C3, NimbleAI will explore the novel concept of *Virtual Neural Networks (VNNs)* to allow users run large and accurate event-driven inference models in only 50 mm² chips. As shown in Fig. 1, this concept will be supported by dedicated TSVs and 3D layers of RAM and NVM that will be architected to create a high-bandwidth and high-density memory hierarchy for quickly swapping active and non-active neurons and (parts of) networks in the dataflow accelerator.

Neuromorphic event- and region-based processing in NimbleAI will help limit the complexity and energy-consumption of AI models, and thus deal with challenges C1 and C2: AI models and algorithms that work on selected image regions are simpler than those that work on full images, and

event-driven networks that execute on neuromorphic hardware only consume energy when there are significant changes in their neuron states, which are themselves triggered only by significant changes in input visual data. Hence, as opposed to state-of-the-practice, which downscale the resolution of input images to keep workloads manageable, NimbleAI will process selected full-resolution image parts for better accuracy. Also, as opposed to state-of-the-practice approaches, where more complex/accurate AI models translate directly into more computing and energy consumption, in NimbleAI model complexity to workload translation will be dynamically regulated through runtime optimization mechanisms that control visual event generation and processing rates along visual pathways.

This unique optimization approach is opposed to the current situation in which performance and accuracy trade-offs are often presented to users as a necessity at the design phase that remains fixed in deployment. NimbleAI will not oblige users to choose between accuracy or efficiency. Instead, it will offer to the user a series of novel runtime system-level optimization strategies that will be continuously refined by means of online learning and applied directly on the user-trained models.

IV. THE NIMBLEAI 3D STACKED ARCHITECTURE

This section describes each of the stacked layers in the 3D NimbleAI chip shown in Fig. 2.

A. Light-field DVS with digital foveation

NimbleAI will implement a digital foveation mechanism to dynamically group and ungroup DVS pixels in the sensor layer to form macro-pixels with varying resolution levels across image regions based on the information value each region brings to the application. If the selective attention algorithms (subsection IV.B) identify something potentially meaningful, DVS macro-pixels in that region will be ungrouped to form a foveated full-resolution ROI that will be processed by a dedicated downstream inference engine. Several foveated regions that match the size, shape, resolution and moving dynamics of the recognized and tracked objects in the scene could be sensed simultaneously to achieve the most accurate results without unnecessarily increasing the amount of data to be processed.

NimbleAI will also be looking at insect compound eye to enable 3D perception for accurate depth and motion estima-

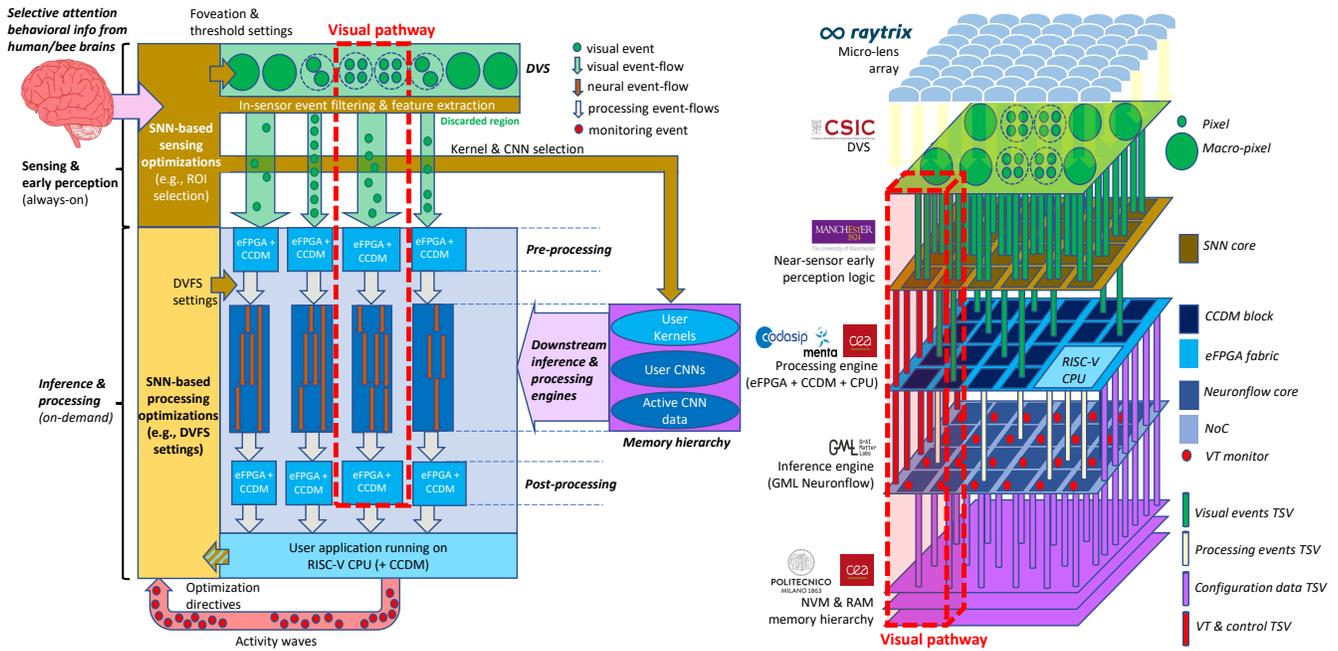


Fig. 2: NimbleAI: conceptual functioning (left) and 3D stacked conceptual architecture (right).

tion. Namely, the project pursues to adapt Raytrix light-field technology (i.e., micro-lenses) [9], and algorithms to encode 3D visual scenes in the form of sparse events that also include depth information: (x,y,z,t) . The fact that DVS events reflect moving edges of the objects in the scene and that light-field algorithms rely on correlations between neighbour data with lots of redundancy, leads us to think that large amounts of processing and energy could be saved by combining both technologies. In fact, it has already been demonstrated that boundaries-first processing lends well to DVS events while putting less pressure on the hardware [11]. Note that depth information opens new opportunities for improving event selection in the early perception and optimization stage (e.g., identify ROIs), as well as for improving perception accuracy in the inference stage.

NimbleAI will investigate 3D silicon integration to vertically stack the additional in-sensor logic needed to implement the functionalities described above, and CEA RRAM [10] to store DVS adjustable parameters (e.g., calibration words and thresholds), with the objective of reducing the impact on the pixel footprint and thus support sensor resolution scaling. NimbleAI will manufacture a DVS testchip with limited resolution and using affordable technology nodes to demonstrate the digital foveation and adaptive region-based sensing mechanisms with adjustable parameters. To demonstrate 3D perception, a light-field-enabled DVS prototype will be manufactured coupling a custom-made array of micro-lenses by Raytrix on a commercial PROPHESSEE sensor.

B. Near-sensor early perception and optimization

NimbleAI will explore uses of SNNs to run ultra-energy-efficient near-sensor visual scene analysis and distill DVS events, namely: (1) remove unwanted events such as noise (e.g., events related to textured bodies); (2) provide initial

feature extraction such as edge detection; and (3) delimit ROIs of almost arbitrary size around collections of identified key features; i.e., selective attention. SNNs are expected to respond rapidly to dynamic changes in the visual input as DVS events will trigger SNN-based processing, helping drive digital foveation in the DVS with extremely low-latency. In fact, it has already been demonstrated that complex visual-cognitive tasks in bees can be modelled with SNN models [12]. Following these findings, NimbleAI will explore models and topologies of comprehensive SNN-like circuits of the bee brain to assess how internally generated oscillations combined with sensory visual events could generate useful types of attentional performance. Modular and scalable SNN topologies will be investigated to extract 3D information from light-field DVS event-flows, analyzing the relations between the size of micro-lenses, DVS pixels (and macro-pixels) and the number of SNN neurons (see section IV.A).

One major objective of NimbleAI is to achieve ultra-energy efficiency by specializing processing to deployment-specific visual events. This involves optimizing (or refining) processing of (new) visual inputs using pre-trained neural networks. In this regard, SNNs are particularly well suited to online training, as their event-based learning rules typically use only information local to the synapse, requiring significantly less computing than the error back-propagation techniques employed to train traditional artificial neural networks [13]. While the focus will be on energy-efficient inference, NimbleAI will also experiment with a range of synaptic plasticity mechanisms such as reinforcement learning and neuromodulation techniques [14], to explore how meaningful visual behaviour can be adjusted on-the-fly, using reward/punishment signals from other parts of the system; e.g., feedback correct/incorrect selection of ROIs from downstream inference engines.

SNNs will also be explored to make inference and processing-related optimization decisions. In fact, the rich temporal dynamics of SNNs are expected to allow them to harness visual continuity in the scene and make more accurate workload evolution predictions. At design time, SNNs will be trained using integrated circuit implementation information, including: (1) energy-performance curves of inference CNNs with different DVFS settings, (2) energy-accuracy curves for different temporal and image resolutions of visual events, and (3) location-dependent thermal dissipation and transmission characteristics in the 3D architecture. At runtime, SNNs will use real-time analytics delivered by activity and V&T monitors embedded in the architecture as feedback for online learning (e.g., slack to deadline when completing processing). NimbleAI will explore mechanisms to capture and represent the state of event-driven neural networks using data-efficient encoding schemes to anticipate activity due to neuron state changes. Furthermore, predictive encoding schemes will also be explored to only transmit errors relative to predictions and thus sustain energy-efficient online learning in SNNs.

Following the same philosophy as for DVS, NimbleAI will design on-chip V&T monitors that generate digital events when they detect temperature and/or voltage variations above or below configurable thresholds. These events will be directly fed into the SNN to take advantage of event-based low-latency processing. Note that SNN-based processing might be especially relevant in next-generation ultra fine-grain DVFS systems to approach brain-like self-regulated energy distribution mechanisms; i.e., anticipate energy needs across regions. Although this concept might have a longer-term impact, we think that providing SNNs with a unified event-based view of both external (i.e., visual) and internal (i.e., voltage and temperature) insights is a very interesting approach to explore holistic optimization decisions that encompass both sensing and processing.

NimbleAI will rely on using SNN software such as NEST/NEURON to carry out model and topology exploration, as well as neuromorphic hardware such as SpiNNaker [14] and commercial spiking-based chips to test the selected models and topologies in real-time applications.

C. Inference and processing

NimbleAI will leverage state-of-the-art event-driven dataflow architectures (i.e., GML NeuronFlow) as main inference downstream engine. As it occurs with SNNs, the type of events that are processed by event-driven dataflow architectures correspond accurately with DVS events, thus maximizing end-to-end efficiency along visual pathways. Recent research has shown that industry standard CNNs designed and trained with popular AI frameworks (e.g., TensorFlow) can be converted to equally accurate event-driven networks with lower computational complexity and hence greater energy-efficiency [15], [16].

Spatial, temporal and neural activation sparsity will be effectively exploited in the event-driven inference engine to improve energy-efficiency and reduce latency. Hence, neuron

compute and event propagation only occurs on sufficiently significant delta events which have not been filtered out because of sparsity. The processing through the dataflow engine proceeds in a systolic array manner, forming “waves” that flow outwards from physical entry points from 3D stacked layers. To support visual pathways and optimally benefit from the DVS foveation approach, the inference engine will be able to run multiple CNNs simultaneously to which visual event-flows are streamed through dedicated physical entry points. At each time, only CNNs that match size, resolution and dynamics of visual stimuli detected in ROIs will be active.

The NimbleAI inference engine will be enabled for running large VNNs with improved accuracy on resource- and area-constrained 50 mm² chips. Enabling VNNs will require to design efficient hardware mechanisms to virtually augment the effective count of neurons integrated in limited chip silicon area. This will be achieved by enabling store and restore network parameters and data on a 3D memory hierarchy that includes stacked high-density RRAM and low-access time RAM layers. The latter memory hierarchy will implement pre-fetching and synchronization mechanisms integrated within the neuron processing pipelines to ensure that neural network parameters and data are accessed and deployed in a timely and efficient manner. Moreover, novel techniques to support compressed synaptic weights, connectivity, and state encoding and storage will be explored to reduce overall data movement. To achieve ultra-high RRAM capacity, NimbleAI will explore and design 3D crosspoint arrays with one-selector/one-resistor (1S1R) memory architectures [17].

Accompanying the inference engine and VNN-supporting memory layers, the NimbleAI architecture will also include a processing layer that implements a RISC-V extensible CPU and an eFPGA fabric for hosting DSP-like pre- and post-processing engines. Besides application-specific processing, this layer will adapt format and properties of incoming visual event flows to best match the available user-trained CNNs and exploit the hardware optimization mechanisms implemented in the inference engine.

The RISC-V CPU and eFPGA fabric will integrate CEA in-memory computing Computational SRAM (CSRAM) [18] blocks to exploit vector computation with less data movement. Furthermore, coupling the CSRAM with eFPGA results in Closely Coupled DSP-Memory (CCDM) blocks provide data parallelism at various granularities to deal with data-intensive operation patterns. Likewise, tightly coupling eFPGA and CCDM with CPU will allow an existing processor design to be specialized for application-specific processing by adding custom instructions (e.g., vector processing and ad-hoc multiply and accumulate) and microarchitecture features, even after deployment. This integrated adaptable processing architecture will reduce data traffic between the CPU and the memory by performing logic, arithmetic, and DSP operations directly in-memory using CCDM. NimbleAI will study the programming model for such an integrated processing architecture that includes CCDM, eFPGA and CPU, and will specify the instruction format generated by the CPU that ultimately define

the user application code. This programming model will be implemented and integrated within the HybroGen software environment for compilation and code generation [19].

D. Physical structure and implementation

A major objective of NimbleAI is to integrate the components explained in previous subsections into an optimized 3D stacked silicon architecture, where each layer is to be implemented using the most appropriate process technology.

To achieve this, the project will develop an EDA tool that supports novel co-design methodologies covering technology-aware 3D architecture exploration across layers and integration to physical implementation. The NimbleAI 3D EDA tool will build upon MZ Technologies Genio 3D tool [20] and third-party physical implementation and analysis EDA tools for 2D IC design. The architecture exploration will consider technology-related aspects, such as process technology and component size trade-offs, and will help make decisions related to layer floor-planning, vertical arrangement of layers, and inter-layer TSV locations to increase computation density and performance, boost communication bandwidth and minimize the length of physical paths. A special focus will be put in designing thermal models of the 3D architecture to pinpoint the locations where to insert V&T monitors to increase the visibility of energy dynamics and thermal dissipation that guide the runtime optimization decisions.

NimbleAI will be looking into integrating IMEC's latest generation of baseline TSV models in a suitable format for the pathfinding engine in the 3D EDA tool to ensure compatibility with silicon technologies, including hardware validation of TSV processes. As shown in Fig. 2, the greatest density of TSVs is expected to connect the DVS sensor with the near-sensor logic layer. A significant less amount of TSVs are expected to connect the near-sensor logic layer to processing and inference engines in the layers below, as well as to support control and monitoring data exchanges among layers. Finally, a medium density of TSVs is expected to support the VNN mechanism, connecting the Neuronflow cores in the inference engine with the memory hierarchy layers. Thermal feasibility and TSV integration limitations will be studied in all these cases.

V. TAKEAWAYS

NimbleAI takes inspiration from ultra-energy-efficient eye-brain systems, even combining divergent evolutionary developments, such as foveation in vertebrate eyes and compound insect eyes. The project expects to deliver 100x energy-efficiency improvement and 50x latency reduction (w.r.t. CPU/GPUs processing frame-based video) by using: (1) DVS sensing with digital foveation and selective attention; (2) event-driven visual inference at optimal DVFS point; (3) specialized processing with in-memory computing; (4) 3D-integrated visual pathways; and (5) system-level optimizations to continuously adjust sensing and processing in each visual pathway to operate jointly at the optimal temporal and data resolution.

NimbleAI will design EDA tools to customize and integrate the technologies and mechanisms above on a sensing-

processing 3D silicon stacked chip. The project will deliver a prototypic implementation of this 3D architecture using an FPGA, small-scale 2D stand-alone testchips and commercial neuromorphic chips. This prototype will be accompanied by the corresponding programming tools to develop and run computer vision applications on it. It will be used as a research vehicle to test novel AI algorithms and runtime optimizations in use-cases related to medical imaging, autonomous driving, eye-tracking and space missions. It is expected that findings coming out from this research will lead to practical implementations in next-generation commercial chips.

ACKNOWLEDGMENTS

NimbleAI has received funding from the EU's Horizon Europe Research and Innovation programme (Grant Agreement 101070679), and by the UK Research and Innovation (UKRI) under the UK government's Horizon Europe funding guarantee (Grant Agreement 10039070). See: <https://www.nimbleai.eu>.

REFERENCES

- [1] M. Tan et al., "EfficientNet: Rethinking Model Scaling for Convolutional Neural Networks," Intl. Conf. on Machine Learning, 2020.
- [2] J. Shalf, "The future of computing beyond Moore's Law," Philosophical Transactions of the Royal Society, 2020.
- [3] N.P. Jouppi et al., "In-Datacenter Performance Analysis of a Tensor Processing Unit," Annual International Symposium on Computer Architecture, 2017.
- [4] C.D. Schuman et al., "Opportunities for Neuromorphic Computing Algorithms and Applications," Nature Computational Science, vol. 2, 2022.
- [5] P. Merolla et al., "A Million Spiking-Neuron Integrated Circuit with a Scalable Communication Network and Interface," SCIENCE, vol. 345, no. 6197, 2014.
- [6] D. Gehrig and D. Scaramuzza, "Are High-Resolution Event Cameras Really Needed?," ArXiv abs/2203.14672, 2022.
- [7] F. Paredes-Vallés et al., "Self-Supervised Learning of Event-Based Optical Flow with Spiking Neural Networks," Intl. Conf. on Neural Information Processing Systems, 2021.
- [8] O. Moreira et al., "NeuronFlow: a Neuromorphic Processor Architecture for Live AI Applications," Conf. on Design, Automation and Test in Europe, 2020.
- [9] Ren Ng et al., "Light Field Photography with a Hand-held Plenoptic Camera," Stanford university, 2005.
- [10] E. Esmanhotto et al., "High-Density 3D Monolithically Integrated Multiple 1T1R Multi-Level-Cell for Neural Networks," IEEE Intl. Electron Devices Meeting (IEDM), 2020.
- [11] C. Kim et al., "Scene Reconstruction from high Spatio-Angular Resolution Light Fields," ACM Transactions on Graphics, vol. 3, no. 4, 2013.
- [12] F. Peng and L. Chittka, "A Simple Computational Model of the Bee Mushroom Body Can Explain Seemingly Complex Forms of Olfactory Learning and Memory," Current Biology, vol. 2, no. 2, 2017.
- [13] J.L. Lobo et al., "Spiking Neural Networks and Online Learning: An overview and perspectives," Neural Networks, vol. 121, 2020.
- [14] Mikaitis et al., "Neuromodulated Synaptic Plasticity on the SpiNNaker Neuromorphic System," Frontiers in Neuroscience, vol. 12, 2018.
- [15] L. Deng et al., "Understanding and Bridging the Gap Between Neuromorphic Computing and Machine Learning," Frontiers in Computational Neuroscience, 2021.
- [16] J. A. Pérez-Carrasco et al., "Mapping from Frame-Driven to Frame-Free Event-Driven Vision Systems by Low-Rate Rate Coding and Coincidence Processing – Application to Feedforward ConvNets," IEEE Transactions on Pattern Analysis and Machine Intelligence, vol. 35, no. 11, 2013.
- [17] A. Fazio, "Advanced Technology and Systems of Cross Point Memory," IEEE Intl. Electron Devices Meeting, 2020.
- [18] J.P. Noel et al., "A 35.6 TOPS/W/mm² 3-Stage Pipelined Computational SRAM With Adjustable Form Factor for Highly Data-Centric Applications," IEEE Solid-State Circuits Letters, vol. 3, 2020.
- [19] <https://github.com/CEA-LIST/HybroGen>
- [20] <https://www.monozukuri.eu/genio-3d>