



Challenges in Mixed-Signal IC Design of CNN Chips in Submicron CMOS*

Angel Rodríguez-Vázquez^{1,2}, Rafael Domínguez-Castro^{1,2} and Servando Espejo^{1,2}

(1) Instituto de Microelectrónica de Sevilla, Centro Nacional de Microelectrónica, Edificio CICA-CNM Avda. Reina Mercedes s/n, Campus Universidad de Sevilla, E-41012 Sevilla (Spain).

(2) Escuela Superior de Ingenieros, Universidad de Sevilla, Edificio Plaza de Armas, E-41012 Sevilla (Spain).

ABSTRACT

Vision machines intend to obtain a representation of their environment through the analysis of a flow of optical images. They are targeted to complete in real-time the processing steps required to pass from the constantly changing sensor data to a semantic description of the scene. This is realized through a sequence of spatio-temporal operations where many nonlinear dynamic interactions are made among the bi-dimensional signal values to: first, determine the local image properties (smoothing, thresholding, edges, motion, colour, texture, motion, etc.); second, obtain the generic scene attributes (boundaries, regions, surfaces, objects, etc.); and, third, build a semantic description from these attributes. It results in a continuous compression of the information, from the raw bi-dimensional input signals to the symbolic representation of the scene.

The contrast observed between the performance of artificial vision machines and "natural" vision system is due to the inherent *parallelism* of the former. In particular, the retina (the "camera" of human vision system) combines image sensing and parallel processing to reduce the amount of data transmitted for subsequent processing by the following stages of the human vision system. Based on this, universities and industries have focused their efforts on the development of new generations of vision artifacts capable to overcome the drawbacks of traditional ones through the incorporation of distributed parallel processing, and by making this processing to act concurrently with the signal acquisition. One possible strategy uses flip-chip bonding of physically isolated sensing and processing devices. Other possibility is to incorporate the sensory and the processing circuitry on the same semiconductor substrate. Silicon retinas, smart-pixel chips and focal plane arrays are members of this class of vision chips. CMOS technologies offer unique features for their development due to the availability of good phototransduction devices and the possibility to realize a large catalogue of linear and nonlinear processing functions by using the different operating regions of the MOS transistor.

Industrial applications demand CMOS vision chips capable of flexible operation, with programmable features and standard interfacing to conventional equipment. The CNN *Universal Machine* (CNN-UM) is a powerful methodological framework for the systematic development of these chips. Basic system-level targets in the design of these chips are to increase the cell density (number of processing cells per unit area), and the cell operation speed. As the technology scales down to submicron all the lateral dimensions decrease by the scaling factor λ , and the vertical dimensions scale as λ^{-a} , where a is typically around 1/2. Consequently, the MOS gate capacitance per unit area increases as λ^a , and the small-signal transconductance per unit channel-ratio increases also as λ^a for fixed bias. As a result one could ideally expect,

$$\text{cell density} \propto \lambda^2 \quad \text{time constant} \propto \lambda^{-2} \quad (1)$$

with constant current and, hence, with no penalty on the power consumption, provided that the voltage ranges remain constants. However, the actual scaling scenario is more pessimistic because of the increased influence of second-order phenomena of small-size MOSTs. Some recently reported submicron CNN CMOS chips feature smaller cell density and operation speed, and larger power consumption, than expected from these formulae. Obviously, there is still room to improve these chips through structural and parametric optimization. However, if precision is a design goal, the cell density and operation speed will be inevitably constrained by mismatch and noise. And these constraints are expected to become harder as the signal dynamic ranges decrease because of the down scaling of supply voltages in deep submicron technologies. This talk will address the challenges involved in the design of CNN chips in submicron technologies.

*This work has been founded by the spanish CICYT under project No. TIC96-1392-C02-02 (SIVA).