

FIG.1

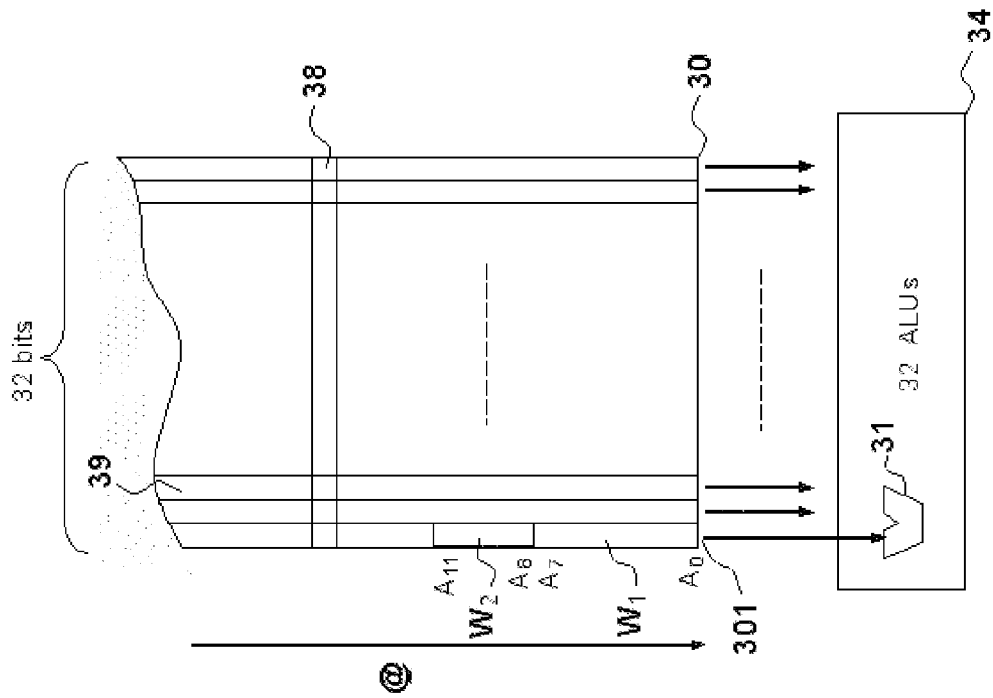


FIG. 2

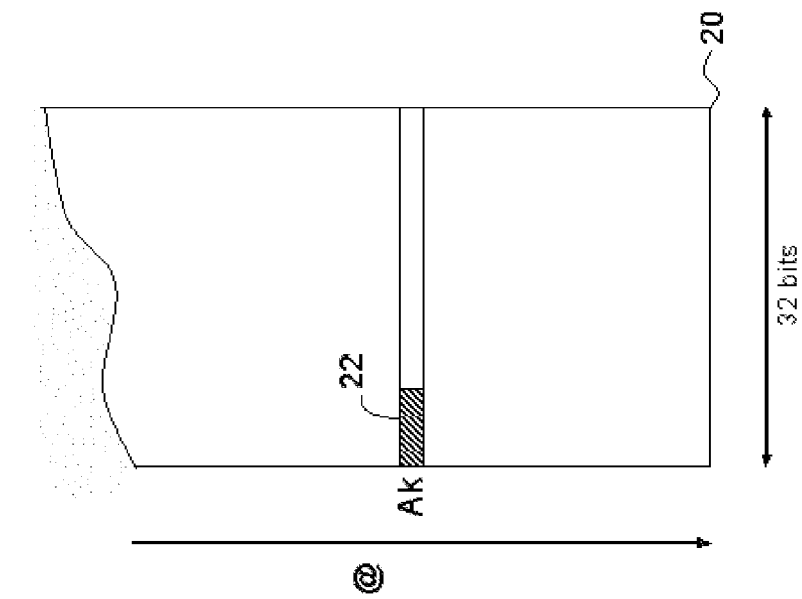


FIG. 3

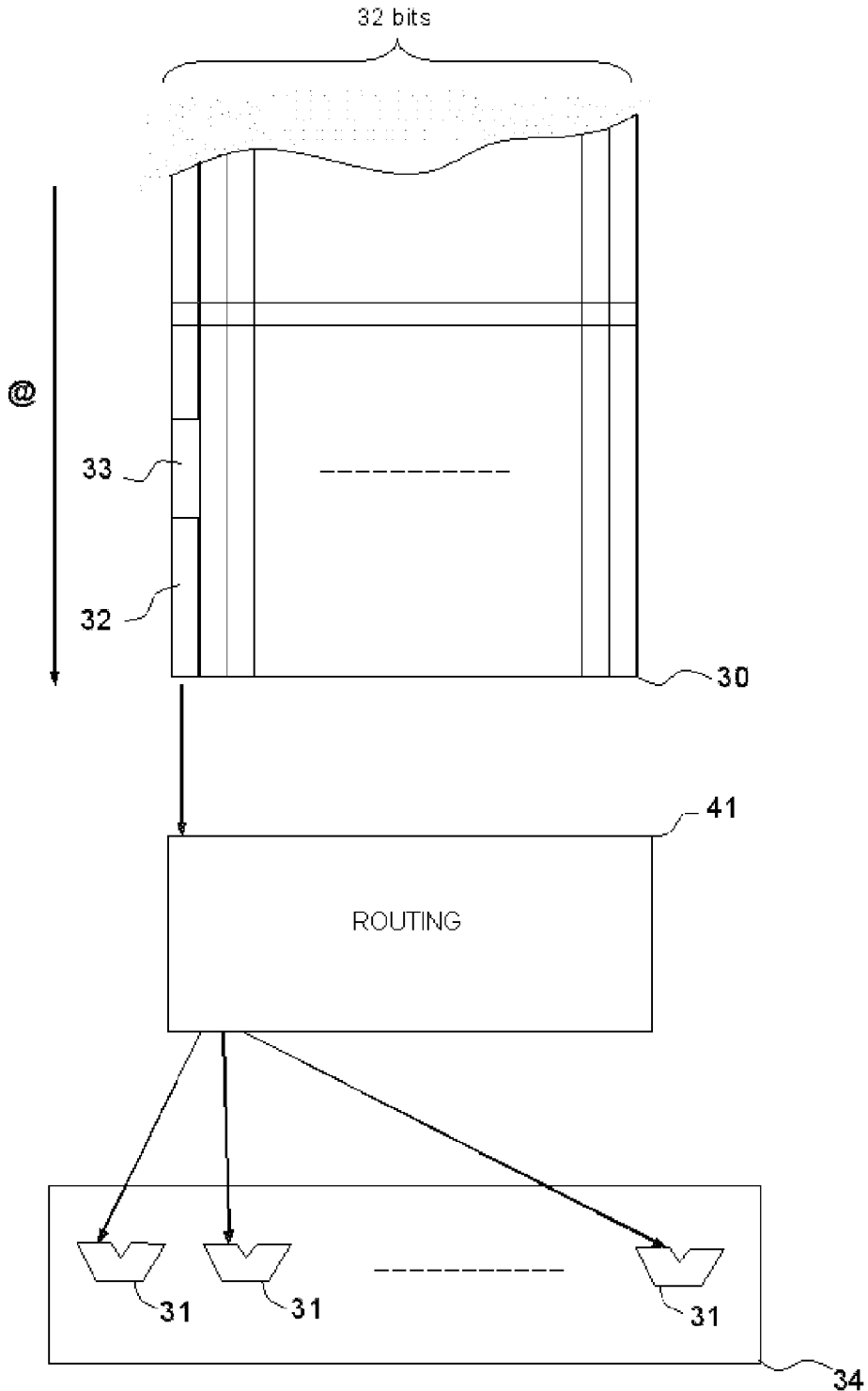


FIG.4

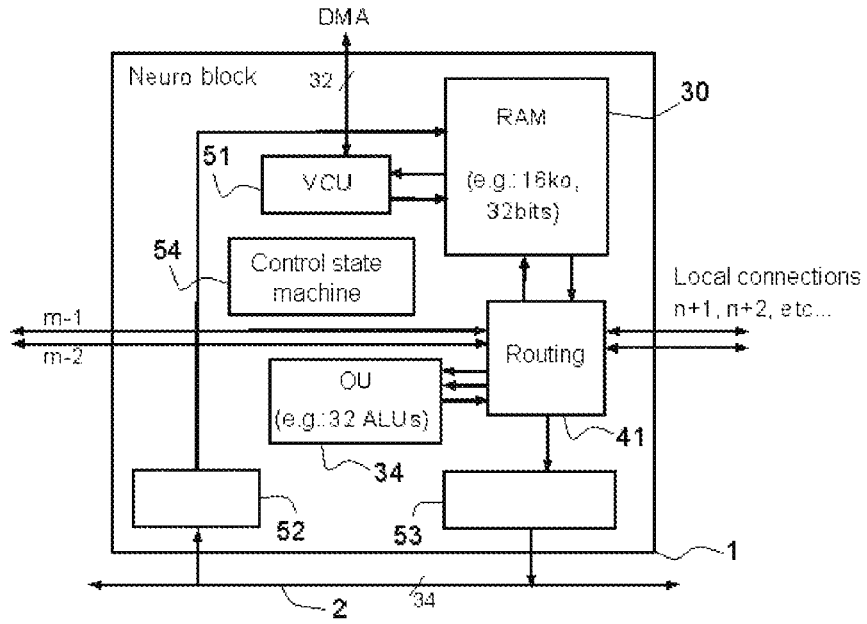


FIG. 5

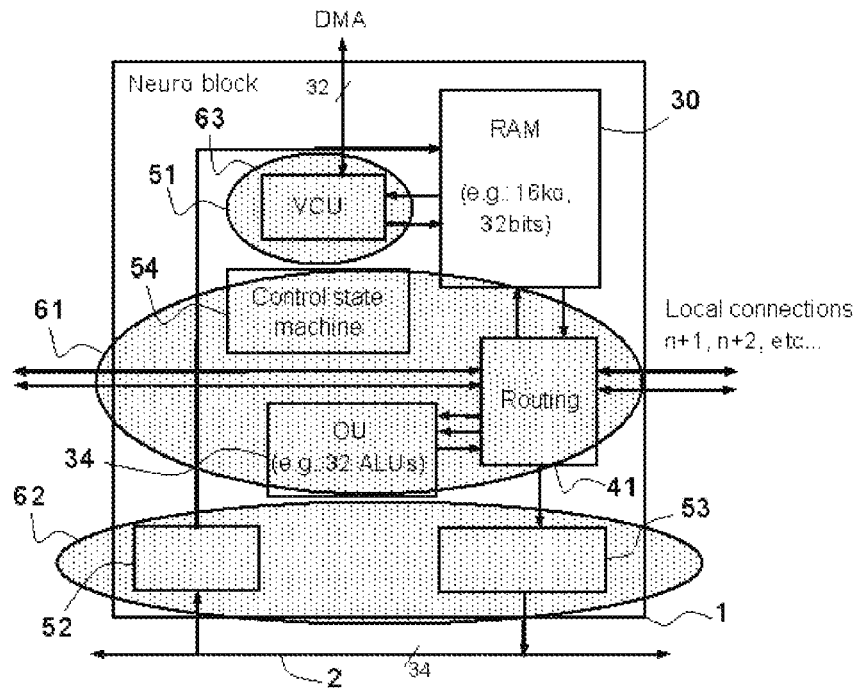


FIG. 6

**SIGNAL PROCESSING MODULE,
ESPECIALLY FOR A NEURAL NETWORK
AND A NEURONAL CIRCUIT**

[0001] The present invention relates to a signal processing module, in particular able to implement algorithms of neural network type. The invention also relates to a neuronal circuit. It applies in particular in respect of the implementation of neural networks on silicon for the processing of diverse signals, including multidimensional signals such as images for example. It also allows the efficient realization of conventional signal processing methods.

[0002] Neural networks are already much used and can potentially be used in very many applications, in particular in all devices, systems or procedures calling upon learning approaches or mechanisms serving to define the function to be carried out, in contradistinction to the more traditional approaches in which the actions to be performed are defined in an explicit manner by a "program". A multitude of systems, ranging from the most sophisticated technical or scientific areas to the areas of everyday life, are thus concerned. All these applications demand ever higher performance in particular in terms of computation power to carry out ever more complex functions, of adaptability, of size and of energy consumption. The algorithms implemented are essential in achieving such performance. The hardware architecture making it possible to implement these algorithms must also be taken into account in the achievement of performance, in particular at a time when the rise in frequency of processors is stagnating or at least seems to have reached its limits.

[0003] To a first approximation, neuronal hardware architectures can be classed according to two aspects:

[0004] A first aspect relates to their structure, the latter possibly being digital or analog, or indeed hybrid;

[0005] A second aspect relates to their specialization in relation to the neural networks liable to be implemented, the architectures possibly being specialized in a few well defined neural networks, such as RBF (Radial-Basis Function) or Kohonen map, or possibly being generic, in particular programmable so as to allow the implementation of a greater variety of networks.

[0006] The types of system addressed by the present patent application are related to generic circuits, with digital implementation.

[0007] The hardware architectures of neuronal systems generally comprise elementary base modules able to implement a set of neurons. In a known manner, a neuron of order i in a neuronal system carries out a function of the type:

$$R_i = f\left(\sum_j w_{ij} E_j\right)$$

[0008] w_{ij} and E_j being respectively the synaptic weights associated with the neuron and its inputs. The elementary module comprises in particular the arithmetic and logic units (ALU) making it possible to carry out all these neuronal functions. f is generally a nonlinear function.

[0009] A technical problem to be solved is in particular that of making efficient use of the silicon on which the neural networks are implanted, especially that of allowing optimal use of the storage of the weights and other data in the

internal memory of the hardware architecture. Another problem is in particular that of allowing a hardware realization that can be expanded as regards the number of neurons/synapses (and therefore of inputs).

[0010] An aim of the invention is in particular to alleviate the aforementioned drawbacks, for this purpose the subject of the invention is a signal processing module, comprising at least one operational unit incorporating computation units, input and output interfaces able to be linked to a bus and a memory storing data destined for said computation units, said memory being organized in such a way that each data word is stored column-wise over several addresses, a column having a width of one bit, the words being transferred in series to said computation units.

[0011] In a possible embodiment, each data word is stored column-wise over several addresses according to an order dependent on the application using said data, this set of several addresses comprises for example address jumps.

[0012] The data transfers are for example performed according to one column per computation unit.

[0013] In a possible embodiment, the module comprises a routing unit connected between said memory and the operational unit, said routing unit having a number of inputs at least equal to the number of width bits of said memory, each input being linked to a single column only, said routing unit routing the data words from said memory to the computation units, one and the same word being able to be routed to several computation units.

[0014] The routing unit comprises for example at least two other series of inputs/outputs able to be linked to circuits outside said module.

[0015] These inputs/outputs are for example able to be linked as inputs and outputs of the routing unit of another module identical to said module.

[0016] The routing unit performs for example all or part of the following operations:

[0017] shifting of the bits of the data words;

[0018] logical operations;

[0019] expansion of the words.

[0020] The module comprises for example a memory virtualization unit linked on the one hand in write and read mode to said memory and on the other hand to an external memory via a circuit of DMA type.

[0021] The memory virtualization unit performs for example operations of reorganization of said memory.

[0022] The operations of reorganization of said memory are for example done by duplication or change of order of the data between the columns of said memory.

[0023] The input and output interfaces communicate for example with said bus by a TDMA protocol.

[0024] Said memory allows for example independent accesses in read and write mode for the input interfaces, the virtualization unit and the routing unit.

[0025] Advantageously, the module operates for example according to several independent synchronous zones, the operational unit operating in a first clock area, the input interface and the output interface operating in a second clock area.

[0026] The memory virtualization unit operates for example in a third independent clock area.

[0027] The computation units execute for example the operations as a function of the value of a guard bit assigned to each of said units.

[0028] The operational unit comprises for example at least 32 computation units, said memory having a width of 32 bits.

[0029] The module being able to implement a set of neurons it performs for example neuronal computations or signal digital processings, said memory storing at least results of the computations, coefficients of filters or of convolution products and synaptic coefficients.

[0030] The subject of the invention is also a circuit able to implement a neural network, characterized in that it comprises at least one series of signal processing modules able to implement a set of neurons such as that described above.

[0031] The signal processing modules are for example grouped together as branches, a branch being formed of a groups of modules and of a dissemination bus, said modules being connected to said bus, a routing block linked to the dissemination buses of said branches performing at least the routing and the dissemination of the input and output data of said circuit to and from said branches.

[0032] Other characteristics and advantages of the invention will become apparent with the aid of the description which follows offered in relation to appended drawings which represent:

[0033] FIG. 1, an exemplary neuronal system comprising a series of elementary processing modules called hereinafter neuro-blocks;

[0034] FIG. 2, the storage of the data (synaptic weights, inputs, etc.) in a memory;

[0035] FIG. 3, a possible mode of operation of a signal processing module according to the invention;

[0036] FIG. 4, another possible mode of operation of a module according to the invention;

[0037] FIG. 5, a possible exemplary embodiment of a module according to the invention;

[0038] FIG. 6, an exemplary manner of operation of a module according to the invention with several independent synchronizations.

[0039] FIG. 1 illustrates by way of example a neuronal system comprising a series of neuro-blocks. The invention is described by way of example in respect of a signal processing module applied to neural networks, but it can apply in respect of other types of processings.

[0040] In the example of FIG. 1, the system 10 comprises 32 neuro-blocks 1. A neuro-block can be considered to be the base element since it is able to implement a set of neurons. As indicated previously, a neuron of order i carries out a function of the type:

$$R_i = f\left(\sum_j w_{ij} E_j\right)$$

w_{ij} and E_j being respectively the synaptic weights associated with the neuron and its inputs and f generally being a nonlinear function.

[0041] In the exemplary layout of FIG. 1, the neuro-blocks 1 are distributed as branches. A branch is composed of several neuro-blocks 1 and of a dissemination bus 2 that are shared by the neuro-blocks linked to this bus. In a configuration with 32 neuro-blocks for example, the neuro-blocks can be distributed as 4 branches of 8-neuro-blocks or as 8 branches of 4 neuro-blocks.

[0042] Moreover, all the neuro-blocks are for example linked by an interconnection line 4 having the structure of a daisy chain bus. More precisely, the arithmetic and logic units (ALU) of each neuro-block can be wired up to this bus. Thus the interconnection line 4 “inter-ALU” passes through all the neuro-blocks 1 of one and the same circuit 10.

[0043] Each branch is linked to a routing block 3, the exchanges between the various branches being done via this block 3. This routing block 3 moreover receives input data and transmits data as circuit output for example via a module for transforming the input/output data 6.

[0044] A direct memory access module 8 (DMA) allows an expansion of the available memory. It is coupled via buses 14, 15 to an internal memory 9, containing a program, and perhaps linked to each neuro-block 1, more particularly to the memory management unit of each neuro-block.

[0045] A control module 11 functions as centralized control processor.

[0046] The exemplary neuronal system of FIG. 1 is used by way of example to illustrate a context of use of neuro-blocks. A processing module, or neuro-block, according to the invention can of course apply in respect of other architectures of neuronal systems.

[0047] FIG. 2 illustrates a problem which arises in respect of a signal processing module, in particular of the neural network type. More particularly, FIG. 2 presents a memory 20 used in a neuro-block. This is a commonplace memory which stores several types of data. In particular it stores the weights of the synapses, coefficients of signal processing filters, in particular carrying out convolution products or Fourier transforms, final or intermediate computation results, as well as other possible data.

[0048] A neuro-block performs very many computations with possibly varying precision. For example in the case of on-line learning, it is possible in particular to distinguish the learning phase, in which the synaptic weights are computed, requiring high precision, for example on 16 bits or 32 bits, and the operational phase which requires lower precision, for example on 8 bits. In any event, the variability of the precision used leads to operations being performed in particular on 4 bits, 8 bits, 16 bits or 32 bits or indeed more and even, conversely, on a single bit.

[0049] FIG. 2 presents a simple example where a word 22 of 4 bits is stored in the memory 20 at a given address, Ak. In this memory having a width of 32 bits, it is seen that the space is poorly occupied. This case illustrates inefficient use of the silicon. By way of example, in a case where the addresses of the memory 20 are coded on 12 bits, the memory has 4096 addresses. For a width of 32 bits, it can then contain 16 kilo-bytes, or more precisely 16384 bytes. In the conventional solutions of the prior art, this available space of 16 kilo-bytes is for example used to 50%, or indeed less if the data are not exact multiples of the width of the memory.

[0050] FIG. 3 illustrates the mode of operation of a neuro-block according to the invention, more particularly FIG. 3 illustrates the mode of storage of the data in the memory 30 of the neuro-block and the mode of transfer to the arithmetic and logic units 31 (ALU). Instead of storing the data by addresses as in the case of FIG. 2, they are stored over several addresses. The example of FIG. 3 illustrates a preferential case where the data are stored on one bit per address, the whole of the data word being stored on several successive addresses. A word W1 of 8 bits is for example

stored on the first bit **301** of the memory between the addresses A_0 and A_7 . Another word W_2 of 4 bits is for example stored between the addresses A_8 and A_{11} . By considering the memory to be a set of rows **38**, each corresponding to an address, and of columns **39**, each having a width of one bit, the memory is filled column by column. Stated otherwise the memory is organized in such a way that each item of data is stored column-wise on several successive addresses from the low-order bit to the high-order bit for example. The filling is thus transposed with respect to a conventional solution where the memory is filled row by row. In the example of FIG. 3, the ranks of the bits increase with the addresses. The reverse is possible, the highest address of the word then containing the low-order bit. Moreover, the words are transferred in series to an operational computation unit **34**, for example consisting of a set of ALUs **31**. A data word is thus tagged according to the rank that it occupies width-wise and the addresses that it occupies column-wise. The word W_2 thus occupies the first bit between the addresses A_8 and A_{11} .

[0051] The transposed filling of the memory, such as described above, combined with the series transfer of the data, makes it possible to optimize the available memory space.

[0052] The storage structure of a module according to the invention, such as illustrated by FIG. 3 affords another advantage. It makes it possible in particular to accelerate certain computations. Because of the "series" type storage inside the memory, it is possible to read a word in either direction. More precisely, the data can be transferred to the computation units starting from the low-order bit, LSB (Least Significant Bit), or starting from the high-order bit, MSB (Most Significant Bit). According to the operations, it is possible to choose a transfer in one direction rather than in the other. Thus for a computation of the maximum between two binary numbers, it is advantageous from the point of view of speed of computation to perform a transfer on the basis of the MSB, the comparisons beginning with the high-order bit. Indeed, if the first item of data has its MSB set to 0 and the second its MSB set to 1, the computation unit can immediately conclude that the second item of data is the larger of the two if the codings are unsigned. On the contrary, for an addition, it is more advantageous to begin the transfer with the LSB, for the propagation of the carry.

[0053] FIG. 3 illustrates an embodiment where the bits are transferred directly to the ALU **31**, the transfers being performed one column per ALU for example. In the case of an application to 32 ALUs with a memory width of 32 bits, one bit rank is assigned to each ALU.

[0054] In the example of FIG. 3, the words are stored on several successive addresses in increasing order. They can of course be stored in decreasing order or disordered if necessary. Moreover, the storage addresses are not necessarily successive, there may indeed be address jumps. In fact, the order and the succession depend in particular on the application.

[0055] FIG. 4 presents a more optimal embodiment where the series words are transferred via a routing unit **41**. This routing unit makes it possible to further improve the use of the memory space. This unit makes it possible in particular to route, or disseminate, a data word toward one or more circuits, in particular toward one or more ALUs **31**. Thus for example, a synaptic weight stored in the memory **30** can be transferred to several ALUs **31** each forming a computation

unit. In particular the filters, convolution products or other types of operations specific to neural networks, have some data in common. Stated otherwise, these data are shared between several operations, for example, one and the same filter is used on several pixels of one and the same image in a convolution operation. In the absence of the routing unit **41**, the aforementioned synaptic weight ought to be stored at several sites in the memory (for example on each column) so as to be transferred to the ALUs which need it. The dissemination effected by the routing unit thus avoids multiple assignments or copies in the memory **30**, a shared item of data being able to be stored at a single site in the memory. **[0056]** Aside from the dissemination function described above, the routing unit **41** can carry out other functions. This unit **41** can for example also perform the following operations:

[0057] Shifting of the bits of the data words (in either direction); for example to facilitate so-called sliding window computations.

[0058] Logical operations;

[0059] Expansion of the words according to various scales by inserting for example one or more "0"s between all the bits of a word.

[0060] The routing unit **41** is for example composed of multiplexers, of registers and logic gates so as to carry out the various data transformation and routing operations. These elements are arranged, in a known manner, in such a way that the operations of routing and transformation, between the memory **30** and the operational unit **34**, can be carried out in a single cycle.

[0061] FIG. 5 presents a possible exemplary embodiment of a signal processing module according to the invention. The module comprises a memory **30**, a routing unit **41** and an operational unit **34** arranged and operating in accordance with the description of FIG. 4.

[0062] The memory **30** is for example a memory of the RAM type having a capacity of 16 kilo-bytes for a width of 32 bits. As indicated previously, it is in particular intended to store diverse data such as input data, results, intermediate results, coefficients of filters or of convolution products for preprocessing as well as synaptic coefficients.

[0063] As in the previous examples, the operational unit **34** comprises 32 ALUs. It carries out all the computations required inside the neuro-block, in particular to perform preprocessings and neuronal processings at one and the same time. The operations which can be implemented are for example the following:

[0064] Addition and subtraction;

[0065] Multiplication and division;

[0066] Calculation of minimum and maximum;

[0067] Numerical computation by coordinate rotation (CORDIC) for trigonometric or hyperbolic functions.

[0068] The ALUs can operate on operands originating from various sources. A first source of data is of course the memory **30**, via the routing unit **41**. The operands can also originate from other modules **1** when the module is implemented in a system, for example a neuronal system of the type of FIG. 1. The transfers between the modules may be done for example by the routing unit **41**. The latter performs for example in the module of rank n a local interconnection with the neighboring modules, for example with the modules of rank $n-2$, $n-1$, $n+1$ and $n+2$.

[0069] The module **1** is moreover able to be connected to an external bus **2**, for example of the type of the dissemi-

nation bus of the neuronal system of FIG. 1, via input and output interfaces 52, 53 which will be described subsequently.

[0070] The neuro-block comprises for example a memory virtualization and neural network topology management unit 51 (VCU) allowing the virtualization of the memory of the neuro-block and the implementation of the various topologies of neural networks. This unit 51 employs direct and independent access to the memory 30, in read and write mode.

[0071] The unit VCU 51 can also ensure global connectivity between the ALUs 34. For this purpose, it possesses a certain number of operators making it possible to reorganize the data stored in memory 30, by duplication or change of order for example (reading of an item of data and writing to another address). It also makes it possible to reorganize the data in memory, for example to replace data which are no longer useful with useful data, allowing for example, the routing unit and the set of ALUs 34 to do the same sequence of operations with the same operand addresses in the memory 30, but with new useful data. The data thus reorganized are ready to be used by the routing unit 41 and the operational unit 34. The unit VCU 51 is moreover linked to a direct memory access (DMA) module outside the neuro-block via a 32-bit bus for example. It can thus read entire blocks in the memory 30 so as to dispatch them to an external memory or write entire blocks to the memory 30 coming from an external memory. The unit VCU 51 thus makes it possible to virtualize the memory containing the synaptic weights, in fact it allows virtualization of the synaptic weights outside the neuro-block.

[0072] The neuro block 1 comprises for example an input module 52 and an output module 53 allowing the connection of the neuro-block to the dissemination bus. In particular, they manage the asynchronism (or in any event the absence of synchronization) between the various modules linked via the dissemination bus 2. In the case of application of FIG. 1, the various modules are, in particular, the other neuro-blocks and, in a silicon implementation, the fact of not forcing the synchronization of the neuro-blocks (that is to say not having an entirely synchronous circuit) will make it possible to gain in terms of operating frequency and to simplify the modular realization independently of the number of neuro-blocks.

[0073] The input module 52 possesses a unique address specific to the neuro-block but possibly reassignable. It monitors in particular the control words of the messages traveling over the dissemination bus: if the neuro-block identifier situated in the header of the message (Neuro-block ID) corresponds to the actual address of the input module 52 or if this identifier corresponds to a dissemination operation, the module captures the whole set of data of the message and stores them in the memory 30 at addresses previously given by the internal program of the neuro-block according to the mode of addressing described in relation to FIG. 3. The output module 53 is equipped for example with a FIFO memory which manages the waiting on account of the TDMA protocol, if the latter is for example used to access the dissemination bus in particular. According to the type of data item, the module 53 may generate a control flit. It locally stores n sequences of 32 bits for example before dispatching them over the dissemination bus, according to the protocol used.

[0074] Advantageously, the TDMA protocol can be combined with the use of the memory 30 by the various resources 51, 41, 52. Indeed, on the one hand the TDMA protocol makes it possible to divide the time into slots, each module internal to the neuro-block having a dedicated slot (for example, a first time slot being reserved for the VCU 51, a second for the routing system 41 linked to the ALUs 34, a third for the input module 52, etc.).

[0075] A block 54, for example an SIMD (Single Instruction Multiple Data) controller with 32 pathways, performs the control of the transfers inside the neuro-block according to a conventional process, known to the person skilled in the art. Moreover, each ALU 31 making up the block 34 is for example controlled by a guard bit, whose state depends on the data to be processed. This guard bit can also be controlled by the block 54. This guard bit allows conditional execution of the operations by the ALUs 31, an ALU executing or not executing an operation dispatched by the block 54 as a function of the value of the guard bit (this guard bit making it possible to disregard the result of the operation if necessary for example).

[0076] FIG. 6 illustrates the various synchronization areas, or clock areas, inside a neuro-block. These various synchronization areas 61, 62, 63 characterize the decoupling of the computations, of the "long distance" communications and of the virtualization. Stated otherwise a first synchronization frequency 61 regulates the computations performed by the operational unit 34, a second synchronization frequency 62 regulates the communications to the dissemination bus via the inputs/output modules 52, 53 and a third synchronization frequency 63 regulates the operation of the memory virtualization unit 51. These three synchronization areas are independent and the synchronization frequencies can vary over time, for example to adapt to the processing speed appropriate to each area in the course of time.

[0077] A module according to the invention allows, in particular, efficient implementation of processings and networks on silicon. The series architecture inside the module and the organization of the storage in the memory 30 allows variable precision to within a bit while optimizing memory occupancy as well as computation time. The invention thus makes it possible to use all the storage resources.

[0078] The decoupling afforded by the various synchronization areas affords the following advantages in particular: increase in the operating frequency of the circuit during implementation, possible variation of the operating frequencies of the neuro-blocks in an independent manner so as to optimize energy consumption, etc. And it makes it possible to decouple the programming of the various parts in the various synchronization areas, thus facilitating the development of applications and the scalability of the architecture proposed according to the various possible realizations (variation of the number of neuro-blocks, of the realization of the communication protocols of the units 52 and 53 etc.).

1. A signal processing module, comprising at least one operational unit incorporating computation units, input and output interfaces able to be linked to a bus and a memory storing data destined for said computation units, wherein said memory is organized in such a way that each data word is stored column-wise over several addresses, a column having a width of one bit, the words being transferred in series to said computation units.

2. The signal processing module as claimed in claim 1, wherein each data word is stored column-wise over several addresses according to an order dependent on the application using said data.

3. The signal processing module as claimed in claim 2, wherein said several addresses comprise address jumps.

4. The signal processing module as claimed in claim 1, wherein the data transfers are performed one column per computation unit.

5. The signal processing module as claimed in claim 1, comprising a routing unit connected between said memory and the operational unit, said routing unit having a number of inputs at least equal to the number of width bits of said memory, each input being linked to a single column only, said routing unit routing the data words from said memory to the computation units, one and the same word being able to be routed to several computation units.

6. The signal processing module as claimed in claim 5, wherein the routing unit comprises at least two other series of inputs/outputs able to be linked to circuits outside said module.

7. The signal processing module as claimed in claim 6, wherein said inputs/outputs are able to be linked as inputs and outputs of the routing unit of another module identical to said module.

8. The signal processing module as claimed in claim 5, wherein the routing unit performs all or part of the following operations:

- shifting of the bits of the data words;
- logical operations;
- expansion of the words.

9. The signal processing module as claimed in claim 1, comprising a memory virtualization unit linked on the one hand in write and read mode to said memory and on the other hand to an external memory via a circuit of DMA type.

10. The signal processing module as claimed in claim 9, wherein the memory virtualization unit performs operations of reorganization of said memory.

11. The signal processing module as claimed in claim 10, wherein the operations of reorganization of said memory are done by duplication or change of order of the data between the columns of said memory.

12. The signal processing module as claimed in claim 1, wherein the input and output interfaces communicate with said bus by a TDMA protocol.

13. The signal processing module as claimed in claim 8, comprising a memory virtualization unit linked on the one hand in write and read mode to said memory and on the other hand to an external memory via a circuit of DMA type, wherein said memory allows independent accesses in read and write mode for the input interfaces, the virtualization unit and the routing unit.

14. The signal processing module as claimed in claim 1, wherein it operates according to several independent synchronous zones, the operational unit operating in a first clock area, the input interface and the output interface operating in a second clock area.

15. The signal processing module as claimed in claim 14, comprising a memory virtualization unit linked on the one hand in write and read mode to said memory and on the other hand to an external memory via a circuit of DMA type, wherein the memory virtualization unit operates in a third independent clock area.

16. The signal processing module as claimed in claim 1, wherein the computation units execute the operations as a function of the value of a guard bit assigned to each of said units.

17. The signal processing module as claimed in claim 1, wherein the operational unit comprises at least 32 computation units, said memory having a width of 32 bits.

18. The signal processing module as claimed in claim 1, wherein being able to implement a set of neurons it performs neuronal computations or signal digital processings, said memory storing at least results of the computations, coefficients of filters or of convolution products and synaptic coefficients.

19. A circuit able to implement a neural network, comprising at least one series of signal processing modules as claimed in claim 18.

20. The circuit as claimed in claim 19, wherein the signal processing modules are grouped together as branches, a branch being formed of a groups of modules and of a dissemination bus, said modules being connected to said bus, a routing block linked to the dissemination buses of said branches performing at least the routing and the dissemination of the input and output data of said circuit to and from said branches.

* * * * *