



UNIVERSITÀ DEGLI STUDI DI TRENTO

DEPARTMENT OF PHYSICS

MASTER DEGREE IN PHYSICS

Optical Bistability As Neural Network
Nonlinear Activation Function

Supervisor:
Paolo Bettotti

Graduant:
Davide Bazzanella

Co-supervisor:
Enver Sangineto

20th March 2018

Contents

Preface	1
1 Artificial Neural Networks	5
1.1 Introduction	5
1.1.1 History	5
1.1.2 Comparison with conventional computers	6
1.2 Basis of Neural Networks	6
1.3 Working Principles of ANNs	8
1.3.1 Learning Process	8
1.3.2 Validation Process	10
1.3.3 Testing Process	10
1.3.4 Datasets	10
1.4 Feedforward NN	10
1.4.1 Other Types of NNs	12
1.5 Real-Life Examples	13
1.6 ANN Simulation	13
1.6.1 PyTorch	13
2 Integrated Photonics	15
2.1 Silicon Photonics	15
2.2 Guided-wave photonics	17
2.2.1 Waveguides	17
2.2.2 Microring Optical Cavity	20
2.2.3 Nonlinear Perturbations	26
2.3 Integrated photonics applied to ANNs	28
2.3.1 Weighted sum of inputs	28
2.3.2 Nonlinear Activation Function	29
2.3.3 Simulations	29
3 Samples, setup and experiments	31
3.1 The samples	31
3.2 Setup	31
3.3 Characterization of the Activation Function	31
3.4 Test of a Trained ANN	32
Conclusions	33
Bibliography	35
Aknowledgements	37

Index

critical angle, 15

effective absorption coefficient, 16

effective refractive index, 16

Free Spectral Range (FSR), 20

Snell's law, 15

Preface

The research field

Research in the field of artificial neural networks has become more and more popular in the last few decades as the computing power available to the average institution increased further. As a matter of fact, in the last two decades the calculation power which once belonged only to bulky, power hungry, and very expensive supercomputers, gradually became possible also for relatively more compact, efficient, and definitely cheaper servers and workstations.

This trend is dictated by the fact that artificial neural networks are usually simulated on conventional computers, which are on the other hand based on the von Neumann, or Princeton, architecture. Implementations of this general purpose architecture are able to do any logical computation with a series of instructions. However the downside of this suitability to generic computations is that it is often difficult to parallelize single operations and therefore certain tasks are inherently inefficient, both in energy and in time. Moreover another problem is that it requires explicit programming for each singular task

The formalization of the von Neumann architecture dates back to 1945 and is as old as the first attempt in artificial neural networks. Shortly after, research and industry efforts focused only on the development of the Princeton architecture, choosing it de facto as the primary design for computing devices.

In the recent years, help came from acceleration units known as Graphic Processing Units (GPU), which were developed for a completely different objective. These GPUs, which are nowadays being designed specifically to accelerate calculations of artificial neural network simulations, are able to carry out a restricted set of instruction compared to CPUs, but are much more efficient both in power and in time, because of their parallel execution.

Today's software implementations

Between the most famous examples of simulated neural network there certainly are the older Watson supercomputer from IBM and the more recent AlphaGo artificial intelligence. Both of them arose to popularity because they defeated us humans to our own games. Specifically the effort of Watson was aimed at prevailing the most strongest human contestants at *Jeopardy!*, an American television game, and succeeded in 2011. Watson was powered by a IBM supercomputer and used 0.22 MW of power. AlphaGO, nevertheless, is a deep neural network that in 2016, through reinforced learning algorithms, defeated Lee Sedol, considered the most formidable master of Go, a complicated board game.

Another interesting case is given by SpiNNaker, of the University of Manchester. It is composed by 10^4 neurons connected by more than 10^7 synapses, and is simulated on over 65 thousand 18-core ARM processors connected together. Unlike the other two, its main goals are to simulate several neural mechanisms, such as the operation of visual cortex. The SpiNNaker system is part of the Human Brain Project (HBP), funded by the European Union.[1]

Similarly to SpiNNaker, another effort in gaining better understanding of the human brain is *The Blue Brain Project*, which is leaded by École Polytechnique Fédérale de Lausanne in collaboration with over one hundred other international research institutions and it is funded by the European Union. This project is simulated on IBM Blue Gene supercomputers, each

fitted with more than 100 thousand processors and consuming several MW of power.

Present objectives and future trends

We can distinguish two principal objectives in this research field. The first one is a technological objective and is that of developing a more powerful computational instrument. The second objective is instead a topic of more fundamental research: to use the tools of artificial neural networks to better understand biological neural networks such as our brain.

If the comprehension of the complex mechanism at the base of our brain has always been a very arduous assignment, the tasks that artificial neural network are facing become increasingly difficult while time passes. Whereas in the past years public and private research groups trained themselves with "easier" problems, such as television or board games, today they are concentrating their efforts on more challenging goals. Probably the most clear example of this tendency is given by the research on autonomous driving, where artificial neural networks are used primarily for real-time processing of data from several sensors.

A great number of businesses formed in recent years around these problems and many other will probably do the same for similar areas. It is very likely that to achieve their objectives, even more powerful and complex artificial networks will be required. Therefore more processors will be developed and put together to achieve impressive computing power. However, one cannot expect to increase indefinitely the power, without significant improvements in efficiency. Hence alternatives are sought.

Overview of current alternatives to simulations

Apart from the vast majority of the research on neural networks, which is being carried out with simulations on powerful processors, a smaller portion of research is focused on building neural network physically. Great efforts are made nowadays by many research institutions and companies to develop original computing architectures that allow artificial neural networks to run physically, instead of being simulated on conventional computers.

The reason why considerable improvements can be expected is that the path of designing hardware devices ad hoc for neural network computation has not yet been covered comprehensively. Major performance or efficiency enhancements are coming from the fact that those new architectures are conceived with parallel execution of specific operations in mind. This, however, makes them inevitably less adaptable to changes, unlike software simulations on conventional computers, which can be modified just by changing few lines of code. A good example is given by size of the network, which in simulation is easily altered, even of orders of magnitude, while in hardware implementations it depends strictly by the architecture features and scalability.

Many of these use the electronics framework, mainly because of the many benefits of the mature silicon technology and its CMOS-compatible production chain. Specifically footprint and efficiency of CMOS devices has not yet stopped to grow since its discovery, in the past century. However, a relatively small portion of these works also inquired simpler network fabrication with organic electronic materials, achieving remarkable results.

Among the most interesting projects which uses electronics there are TrueNorth of IBM, Neurogrid of Stanford University, and BrainScaleS designed in Heidelberg and Dresden. The former is part of a long-term research project funded by DARPA and has the aim of building an artificial neural network with the capabilities of brains of small mammals, such as cats or mice. In 2014 IBM showed a chip, named TrueNorth, containing 1 million artificial neurons and 256 millions synapses. This network, relying on more than 5 billion transistors organized on an area of 4.3 cm^2 , consumed only 60 mW. The researcher also proved that connecting 48 of this chips together they obtained the equivalent of the brain of a mouse.

Similarly, the intention behind Neurogrid chip is to explore numerous hypothesis concerning mammalian brains, specifically about the inner mechanisms of operation of the cerebral cortex. Stanford's implementation reaches the line of 1 million neurons, while only requiring 5 W of power consumption. Moreover the communication lines between single neurons, the synapses, are provided by FPGAs and banks of SRAM.

Finally the BrainScaleS is built by 20 silicon wafers, each containing 200,000 biologically realistic neurons and 50×10^6 plastic synapses. Likewise SpiNNaker is funded by the European Union and it is part of the Human Brain Project. A particularity of this system is that does not execute pre-programmed code but evolves according to the physical properties of the devices.[1]

Besides electronics, other kind of physical implementations are attempted from research teams. Probably the most intriguing example is given by researches in photonics, which attempts to match electronics performance by overcoming its intrinsically weaknesses, such as power consumption, speed, and intrinsic compatibility with parallel computing. Photonics has yet to reach achievements comparable to electronics, mainly due to the embryonic state of its technology.

My work

My work in this thesis has two related objectives. The primary aim is to implement a proof of concept for a fundamental component of artificial neural networks: the neuron's activation function. I plan to do so by employing the framework made available by integrated photonics. The second intent is to bring together two until now distant fields of research: physics and information technology, specifically photonics and machine learning. In fact generating knowledge transversal to the two research fields is at least as much important as the primary objective.

The idea at the base of our physical implementation is to employ integrated optical devices, which have been manufactured for other projects, to create a proof of concept for the activation function. Therefore in this work I will study a particular device, a microring resonator, and characterize its response to different inputs, for several working conditions. This device can process information because it has, under certain circumstances, a nonlinear response function to its inputs. Moreover, the same device might be used in another important part of artificial neural network nodes: the weighted sum. With this two parts together one could claim to have built a rudimentary artificial neuromorphic network, the *perceptron* (see 1.4). However, since integrated structures able to carry out a weighted sum of optical signals have already been discussed in literature and it requires less effort to produce acceptable results, it will not be at the center of this study.

During the thesis I achieved two main results. I have demonstrated the use of optical bistability in microring cavities to produce an activation function I have demonstrated the use of the response of a microring optical cavity to produce an activation function by exploiting the optical bistability regime. In addition, with the aim of testing such response function, I simulated an artificial network with custom activation function by using standard software libraries used by ICT community to model neural networks.

Structure of the thesis

The first chapter of this thesis is intended as an introduction to the world of artificial neural networks. After a brief historical summary, I describe the main aspects of neural networks and their working principles. Then I focus on the description of the simplest type of artificial networks, the feedforward network, and finally and overview on other kind of networks. Later, the chapter is concluded with an introductory description of the language Python and the library PyTorch, used in this work for artificial network simulations.

The second chapter introduces the physics which I will study. Again, after a short explanation on the basic devices used in the field, I will describe in depth the device studied for the task. Having clarified the physics underneath, I ought to explain how I am going to use this physics to implement (part of) an artificial neural network. I will address these topics both from the theoretical point of view and from the numerical one, with simulation of the specific (simplified) system.

The third chapter is dedicated to the description of the experimental work. At the beginning the problem of the selection of the sample is discussed, along with the definition of the device chosen. After, I characterize the different behaviors with detailed data. Finally, an explanation on how the device will be used in the neural network viewpoint and the corresponding tests of operation.

Chapter 1

Artificial Neural Networks

1.1 Introduction

Artificial Neuromorphic Networks (ANNs) are computational systems which elaborate information in a way that is loosely inspired by the operation of biological neural networks (animal brains). Biological networks are superior in performance to computers and extremely more efficient in difficult tasks such as classification (e.g. image recognition) and prediction (e.g. pattern recognition). The underlying idea is to copy some of their mechanisms and exploit them for computational applications. These systems are intrinsically parallel in operation, do not require specific programming to operate, and modify their behavior through a learning process to improve their accuracy in a certain task.

Mathematically speaking, ANNs are a collection of nodes, each one of them elaborates the information and is somehow connected to the other. Artificial networks can be either simulated on computers or physically built on hardware designed ad hoc. At this time, ANNs are mainly implemented by simulations, however the technological progress is withheld by limitation in computing power and efficiency. Training a complex artificial neural network with computers at the state of art might take even weeks. Thus, with the aim of improving performance from simulated networks, research on hardware architectures is surging: digital, analog, electrical, and optical devices are being developed.

1.1.1 History

It is widely acknowledged that the opening work of this research field was made in 1943 by Warren McCulloch and Walter Pitts, a neurophysiologist and a mathematician respectively. In their research work, they described the operating mechanism of biological neurons by modeling a simple electronic circuit [2].

The following important work was made by Donald O. Hebb, who hypothesized the neural plasticity in 1949 [3]. He pointed out the fact that neural pathways are strengthened each time they are used, a concept known as *Hebbian learning*.

During the late 1950s and the early 1960s, as computers became more powerful, many promising works were published. Some simulated operation of artificial networks on calculators, e.g. Farley and Clark [4] and Rochester [5]. Other produced circuitry that implemented on hardware such networks, e.g. Rosenblatt [6], [7]. However, despite the early successes of neural networks, the traditional computing architecture (von Neumann architecture) was chosen as the preferred computing architecture.

The reason why this happened, probably, was due to several concurring facts. First of all, in the same time period (1969) a research paper by Minsky and Papert [8], which identified two important problems. The basic perceptron was not able to execute the exclusive-or (XOR)

operation, unlike the logical circuits at the base of von Neumann architecture. Moreover, the research stated the fact that more complex networks such as multiple- or deep-layered networks were not possible (at that time) because of the lack of adequate processing power.

In addition to this research, the early successes of some works on neural networks pointed to an overestimation of the artificial neural networks potential, also held back by the technological capacity of the time. Finally, important questions of philosophical nature came to light, such as the fear fueled debate on the impact on our society of a class of computers able to think. This very controversy, i.e. the artificial intelligence (AI) problem, is discussed still today [9].

Sometime during the 1980s the interest for this computing method was reinvigorated. The main stimulation was probably given by a number of works which suggested methods to implement multi-layered networks distributing pattern recognition errors through the all the layers in the network. This method is now called *backpropagation*.

In today's research and technology, artificial neural networks are used in numerous applications. However, the development is made slowly, due to technological limitation in computational power of present processors.

1.1.2 Comparison with conventional computers

1.2 Basis of Neural Networks

A neural network is a collection of processing elements, or nodes, interconnected in an arbitrary topology. From its input nodes, the network accepts information, which will propagate into the inner nodes through the interconnections and will get elaborated at each node. At the end of the network, there will be a number of output nodes, with the task of reading a portion of the inner nodes. The inner nodes are also called hidden, because they are not meant to be accessible to the external world. A generic scheme of such network is shown in Figure 1.1 on this page.

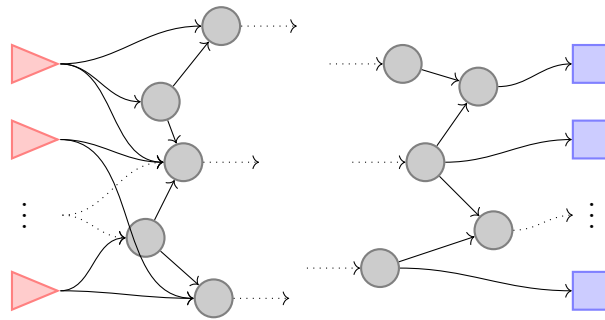


Figure 1.1: Generic scheme of a neural network. Triangles (red) are input nodes, circles (grey) are inner nodes, and squares (blue) are output nodes. Interconnections among nodes are represented by arrows: continuous when both elements are drawn, and dotted otherwise.

Nodes can all implement the same function or behave differently, depending on the type of neural network. The operation of nodes resembles that of animal neurons: various input gets collected and elaborated together to obtain an output, which will become one of the many inputs for subsequent neurons/nodes. Specifically, the most used model for neurons is the McCulloch–Pitts (MCP) neuron. It is divided into two parts, as shown in Figure 1.2: the first part is a weighted sum of the inputs, while the second part is given by the so called activation function.

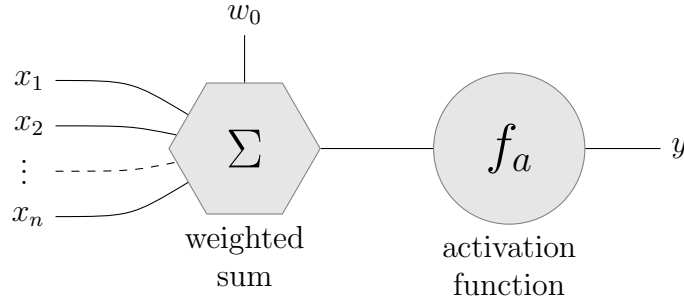


Figure 1.2: Generic node representation. x -values are inputs, y -values are outputs, w_0 is the bias.

The node is described mathematically by equation (1.1)

$$y = f_a \left(w_0 + \sum_{i=1}^n w_i x_i \right), \quad (1.1)$$

where f_a is the activations function, evaluated on the sum of the input x_i weighted with w_i , plus a bias w_0 .

Each node accepts values at its inputs and produces an output accordingly. However, in addition to the input, the output depends also on the node's parameters: the weights and the bias, which are usually changed outside the operative phase of the neural network (see Section 1.3).

Moreover it is mandatory for the activation function $f_a(\cdot)$ to be nonlinear, because otherwise a collection of nodes will result in just a weighted sum of its inputs. Two examples of nonlinear function are shown below in Figure 1.3.

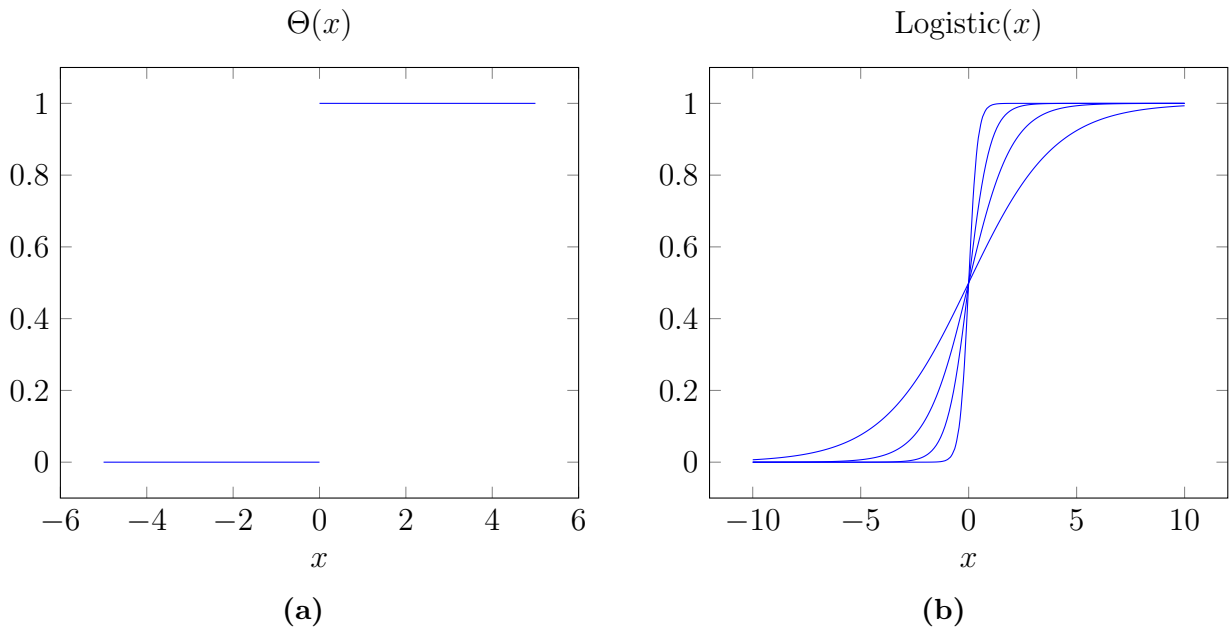


Figure 1.3: Examples of activation function: (1.3a) is the well-known step function, or Heaviside Θ , (1.3b) depicts a few functions from the family of the Logistic functions.

One can distinguish at least three type of nodes in every neural network: input, inner/hidden, and output nodes. Input nodes take one input value, from the outside of the neural network, and pass it on to the inner nodes unchanged. Inner/hidden nodes take many

inputs and generate an output through the activation function. Output nodes, similarly to input nodes, take one input value, from the inside of the NN, and pass it on to the outside.

Standard Representation

The way I depicted a generic neuromorphic network in Figure 1.1 is not the standard representation used in books and research papers. The main difference is that usually weights are commonly represented on the connection between the nodes, which are then designated to apply only the activation function. Moreover the input layer is linear as it feeds the inner nodes with the input data, while the output layer is actually given by the last nonlinear layer of the inner nodes. A generic network is shown in Figure 1.4a.

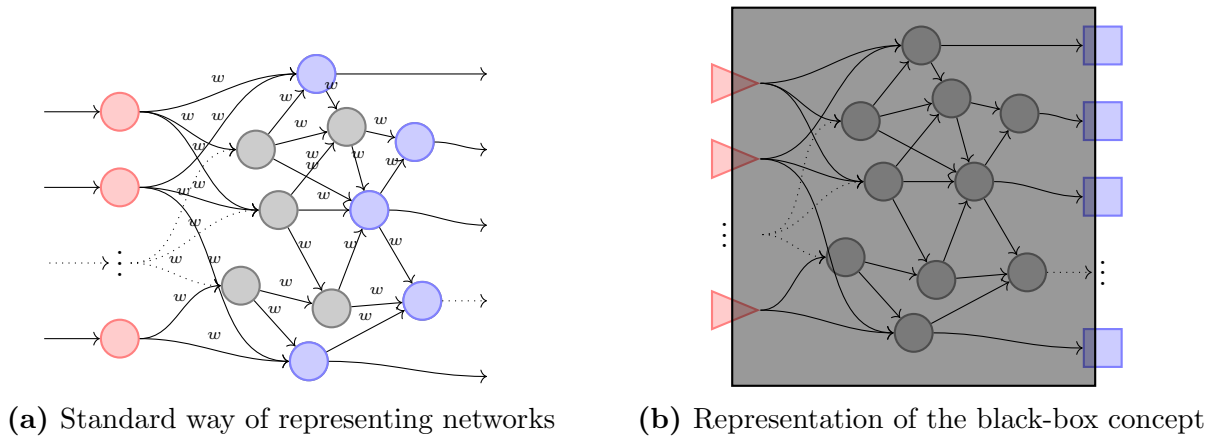


Figure 1.4

On the contrary, I consider the inner nodes as the only place where any kind of elaboration on the data happens. Inner nodes have a number of inputs, which are weighted and summed together to be entered as argument in the activation function. This leads to a natural separation between input/output nodes, which acquire the task of providing data from/to the outside, and inner nodes, which is where the activation function and/or the weighted sum are carried out.

This non-standard description, moreover, is consistent with the idea of functional *black box*, in which input and output are the only visible nodes, while the other are hidden inside, as shown in Figure 1.1.

1.3 Working Principles of ANNs

Because of its topology, each neural network will behave in a different manner from other neural networks with diverse, or even similar, arrangements of nodes. Moreover the same neural network will perform a certain task better or worse also depending on how inputs are weighted at each hidden node, and normally those parameters are initialized with a random value at the creation of the network. For this reason, before a neural network is considered ready to perform a task, it usually must go through three training stages: learning phase, validation phase, and testing phase. Every one of these stages is meant to prepare the network to work as required from the designer.

1.3.1 Learning Process

During the learning process the neural network is run on a set of known inputs x , each paired with its correct answer y , or target, in a second set of data. The neural network will produce

at the output a third set \hat{y} , which should be as close as possible to the correct answers, when the network works properly. Typically weights among nodes are initialized randomly and the distance of the network outcome from the target function is measured through a loss function. Then weights are modified in order to decrease the loss function to its lowest possible value.

Loss function

The loss function $L(y, \hat{y})$ evaluates the difference between the predicted and the correct answer. Usually, this quantity is linked to the geometrical distance between the predicted output and the target $|\hat{y} - y|$.

The most common loss function is the *mean-square error*. Assuming to have an input set of N examples paired with the same number of targets, and that the outputs and the targets are composed by C values, or classes, the function becomes:

$$L(y, \hat{y}) = f_{MSE}(y, \hat{y}) = \frac{1}{N} \sum_{n=1}^N \sum_{i=1}^C (\hat{y}_{n,i} - y_{n,i})^2, \quad (1.2)$$

where each example in the set is subtracted to its target and then squared. Finally the mean of all squares gives the expected result.

Another commonly used function is the cross-entropy loss (also known as negative log likelihood),

$$L(y, \hat{y}) = f_{CEL}(y, \hat{y}) = -\frac{1}{N} \sum_{n=1}^N \sum_{i=1}^C y_{n,i} \log(\hat{y}_{n,i}), \quad (1.3)$$

which expects positive values at the input. Hence the error $-y \log(\hat{y})$, quantified for each element in each example, is always a positive number. The mean over all examples in the set returns the results.

Alternatively, variations of the previous methods are given by taking the sum of the examples in place of the mean, or by calculating a loss function for each example instead of evaluating it for the whole set.

Weights Update Process

The weights update process is a difficult task, and probably the most computationally expensive one in running a neural network. There is a variety of methods to chose from, depending on the type of artificial network and the resources available.

A widely used algorithm is the gradient descent, from its most simple version to more complex variations such as stochastic gradient descent (SGD). This method updates the weights by subtracting a value proportional to the gradient of the loss function in respect to the weights themselves times a positive factor called *learning rate*, as shown below.

$$w_i|_{n+1} = w_i|_n - lr \cdot \frac{\partial L}{\partial w_i|_n} \quad (1.4)$$

where $w_i|_n$ are the current weights, lr is the learning rate, $\frac{\partial L}{\partial w_i|_n}$ is the first derivative of the loss function in respect to the i -th weight at the current step, and $w_i|_{n+1}$ are the updated weights. This method is equivalent to minimize the error on the loss function, by following the gradient $\nabla_w L$. This vector lives in the multidimensional space of the loss function $L : \mathbb{R}^W \mapsto \mathbb{R}$, where W is the total number of parameters in the network.

The most efficient ?? algorithm is called *backpropagation*: it computes the first derivative of the loss function L in respect to all the parameters of the network, the weights, starting

from the end of the artificial network and going backward toward the input, hence the name backpropagation. Since the number of connections between nodes might be even order of magnitude bigger than the number of nodes, it is simple to understand how large networks are computationally expensive to train.

Other types of learning processes are used, e.g. unsupervised/reinforced.

1.3.2 Validation Process

1.3.3 Testing Process

At the end, there is the process of testing the artificial neural network. Ideally the network is tested on a new set of data, for which the target results are known, similarly to the preceding phases. This time, however, the predicted outputs are compared to the correct answers to obtain an overall value for the correctness, often expressed in percentage.

1.3.4 Datasets

Since there are three phases of preparation for any artificial neural network, there must be an appropriate number of examples to feed to it.

HOW DO I DIVIDE A DATASET?

WHAT HAPPENS WHEN THERE ARE TOO FEW EXAMPLES?

AND WHEN THERE ARE TOO MANY?

1.4 Feedforward NN

The first and most simple type of neural network is called Feedforward. In this kind of neural network, nodes are divided into groups called *layers*. A layer is a collection of nodes that accepts inputs from a preceding group and generate as many outputs as the number of nodes in the layer. Each layer of a Feedforward neural network is connected in series with the others, except of input layer at the beginning and the output layer at the end. As for the single nodes, the inner layer are called hidden, because usually not accessible.

The information travels from the input to the output and gets elaborated from each hidden layer: there are neither connection between nodes of the same layer, nor loops or feedback between layers. The number of hidden layers and the number of nodes they contain depends on the network topology. Moreover the connection between the layers might be complete, i.e. each node in the layer accepts each input of the preceding layer, in that case the layer is said to be *fully connected*, or sparse as in the case of convolutional layers (see Section 1.4).

Perceptron

The most naive topology of a Feedforward neural network is given by the so called *Perceptron*. The Perceptron dates back to the 1957, when the homonym *Perceptron algorithm* was software implemented by Frank Rosenblatt on a computer (IBM 704) and only subsequently in hardware as the *Mark 1 perceptron* [6], [7]. The graph of a generic (single layer) perceptron ANN is shown in Figure 1.5 below.

By adding more than one hidden perceptron layer to the neural network, one obtain the so called *Multi-Layer Perceptron* (MLP). This allows for more computational complexity ?? When the total number of layer is more than two, the network is called *deep*. A deep MLP is shown in Figure 1.6. In principle any shape is possible, i.e. each layer could have a different

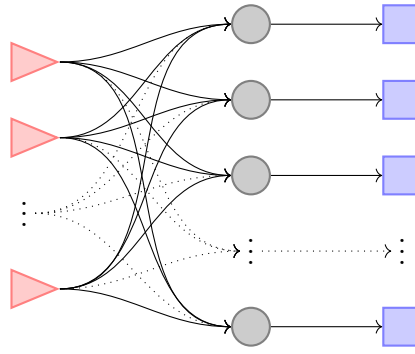


Figure 1.5: Perceptron type neural network: in this representation the perceptron has n inputs and m outputs as well as a hidden layer with m nodes.

number of nodes, however often the layers at the beginning are wider than the layer at the end of the network ?? Besides the shape, in literature a perceptron is almost always considered fully connected ??

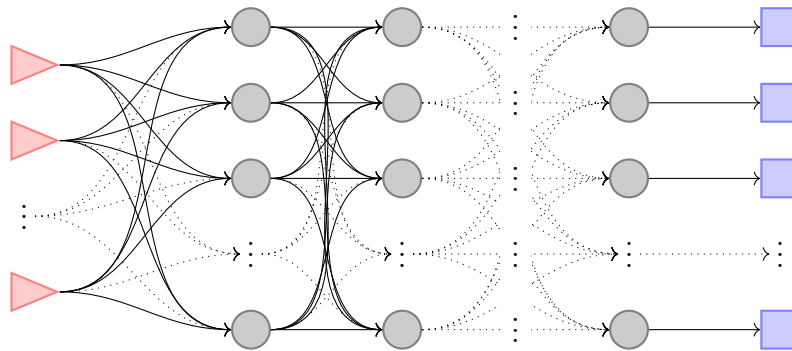


Figure 1.6: Deep Multi-Layer Perceptron (MLP), fully connected.

Other Feedforward NNs

Autoencoder neural networks are feedforward networks in which the output nodes are as many as the input nodes. The purpose of this kind of network is to reconstruct its own inputs.

Probabilistic

Time delay

Convolutional networks are inspired to the visual cortex, in which neurons are not fully connected their inputs but only to a restricted region. Convolutional neural networks are a type of feedforward network conceived to recognize images without being misled by distortions such as translation, skewing, or scaling. In this kind of network the input is often represented by a 2D matrix, instead of a 1D vector. It is usually composed by many layers, the prevalent kind is the convolutional one. This layer performs a two-dimensional convolution over the input matrix of a second 2D matrix of weights, called *feature map*. Thus, each node of the layer operates on a restricted region to understand if a feature is present or not. Commonly the operating regions are overlapping and the feature map is shared among the nodes in the same layer.

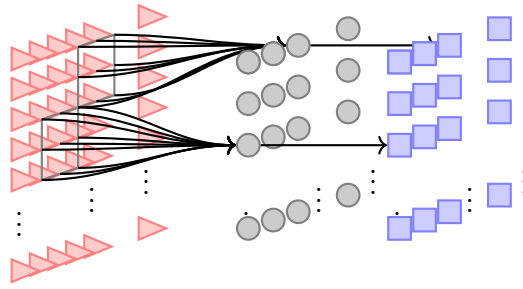


Figure 1.7: Pictorial representation of a layer of a convolutional supernode. Several supernodes form a convolutional layer.

1.4.1 Other Types of NNs

By changing the topology of the nodes distribution and their connections, one obtain other networks that cannot be catalogued under the class of feedforward networks. Moreover, those different types of network are not a niche, but are widely ?? studied as a different approaches to the same or additional problems.

Recurrent NN

Recurrent neural network are a kind of network in which a portion of the input of nodes depends on the (past) output of the same nodes or nodes of subsequent layers. That is information does not propagates only forward like in the feedforward networks, but can propagate also backward, for example in loops or in feedbacks.

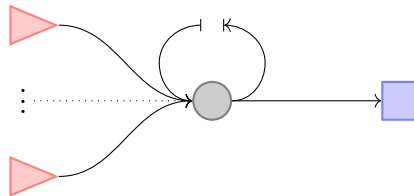


Figure 1.8: Representation of a recurrent node. One of the inputs is given by the output itself. This output to input connection could be mediated by a delay device, so that for example the output at $t - 1$ becomes the input at t . Depending on the structure of the network, there might be recurrent nodes and/or recurrent group of nodes, loops.

Reservoir NN

Reservoir neural networks differ from feedforward and recurrent networks in the learning approach. In fact, the topology of a reservoir network could be exactly the same as that of a deep multi-layer perceptron or that of a recurrent network. However, the reservoir computing differs in approach in respect to deep learning. It claims that it is not necessary to learn all the weights of the network, as in deep learning, but it is sufficient to train only the last (perceptron) layer of the network.

This kind of networks, then, can be trained much faster than their respective counterparts, i.e. feedforward and recurrent. The question over which training method is correct is still debated and literature does not provide clear answers yet.

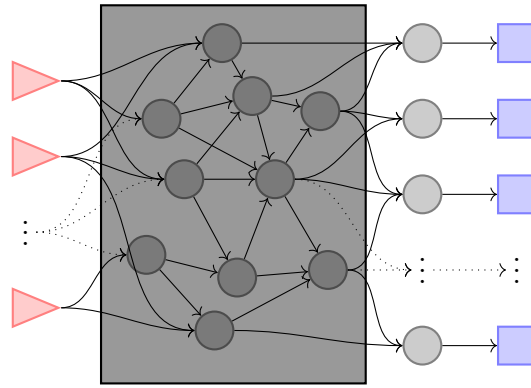


Figure 1.9: A reservoir NN is topologically equivalent to deep networks. However only the last layer is trained. In this picture the trained node are represented outside the box, while those inside are initialized with random weights which are then left unchanged.

Modular NN

Spiking NN

Spiking artificial networks are the most different kind in respect to all the other networks until now described. In this class of ANNs, information is not coded only in the intensity of the signal, but also in the rate of signals, e.g. a high value will be encoded as a signal with high repetition rate, whereas a low value as a signal with low repetition rate. This way of encoding information is more alike the mechanism of biological neural networks, such as our brain ??

LOOK INTO other type of neuron model: Hodgkin–Huxley (H-H) model

https://en.wikipedia.org/wiki/Binding_neuron

1.5 Real-Life Examples

? SHOULD I KEEP THIS SECTION ?

1.6 ANN Simulation

The resources and the time available allowed me to implement physically only the activation function of one node. Hence, to test this hardware implementation as I will show in Section 3.4 of Chapter 3, I had first to simulate and train *offline* a specific neural network. To do so, I chose a programming language, *Python*, and a library, *PyTorch*, which helped me in this task.

1.6.1 PyTorch

PyTorch is a python package that provides high-level features such as a deep learning research platform [10]. It is based on a backpropagation algorithm called *Reverse-mode auto-differentiation*, which is fast and flexible. Moreover integrates acceleration libraries that allow fast and lean operation.

To summarize, the library was chosen for its language and its flexibility, even though it is also powerful. This because our need was to simulate a little network to test the activation function.

Chapter 2

Integrated Photonics

To build an artificial neuromorphic network one has to choose first some physical phenomenon to employ in the fundamental blocks, likewise transistors in electronic circuits use the various behaviors of electrons in resistances, capacitors, and inductances. The physics which I want to build my artificial neuromorphic network with is photonics, precisely integrated silicon photonics.

Photonics is the physical science which studies detection, manipulation, and emission of light. Specifically, integrated photonics is the branch that studies how to reduce and *integrate* once macroscopic optical devices in miniaturized structures. In the past few decades, thanks to the constant improvement of the manufacturing techniques, many productive problems has been progressively resolved. Moreover interest in the field is rising, driven by the growing needs of communication technology, which was in turn following the increase in computational power of electronics. Many integrated devices have been proposed ?? and some of them even commercialized ?? .

Silicon photonics is the research field that studies how to integrate optical devices in structures built with silicon and materials derived from it. Silicon is well-known material which has been widely studied in microelectronics. It possesses many qualities which allow relatively easy manufacturing of high grade structures. However there also some drawbacks, for example being a centrosymmetric element it does not possess nonlinearities of the second order.

2.1 Silicon Photonics

Silicon-on-Insulator (SOI) photonics is one of the widest branches of integrated photonics. It pursues the same objectives of all integrated photonics with the addition of following the production steps already developed for microelectronics.

Its framework is the silicon wafer, which is then altered by a often numerous series of production steps. The first step is the burial of an insulator layer of SiO_2 , called *buffer*, which defines the bottom surface of the layer at the top. The pure *Si* volume at the top, the *device* layer, is where integrated optical structures will be built with various processes like photolithography, thermal oxidation, ion implantation, and etching.

The fact of relying on the SOI framework to build integrated devices is a strength of silicon photonics, because it enables researchers to develop CMOS compatible integrated optical structures. Complementary Metal-Oxide-Semiconductor (CMOS) compatibility is an industry standard created to fabricate microelectronic devices. Therefore by exploiting this standard all the know-how of the manufacturing industries behind the commercial microelectronic products is available to fabricate integrated photonics devices.

Moreover, by sharing the fabrication technology with microelectronic makes it easier to integrate electronic components to obtain hybrid integrated optoelectronic devices or to compensate

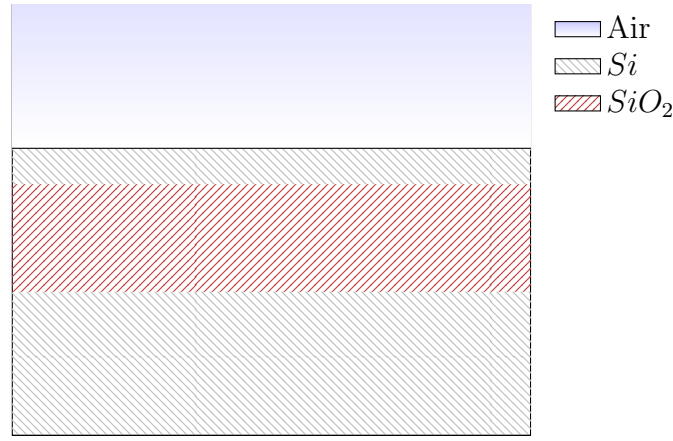


Figure 2.1: SOI productive framework. The buffer layer of SiO_2 is buried underneath the pure crystalline Si device layer which is used to build structures. The top silicon layer will become the volume in which light will be confined. The bottom silicon, underneath the buffer layer, is called substrate.

lack of appropriate optical structures.

Hence silicon photonics is a promising technology for a combination of cost and technological reasons. A list of advantages and disadvantages of silicon photonics in compared to other photonics technologies is provided in Table 2.1.

Disadvantages	Advantages
<ul style="list-style-type: none"> i. Stable, well-understood material ii. Stable native oxide available for cladding and electrical isolation iii. Relatively low-cost substrates iv. Optically transparent at important wavelengths of $1.3\ \mu\text{m}$ and $1.55\ \mu\text{m}$ v. Well-characterized processing vi. Highly confining optical technology vii. High refractive index means short devices viii. Micromachining means V-grooves and an effective hybrid technology are possible ix. Semiconductor material offers the potential of optical and electronic integration x. High thermal conductivity means tolerance to high-power devices or to high packing density xi. Carrier injection means optical modulation is possible xii. Thermo-optic effect means a second possibility for optical modulation exists 	<ul style="list-style-type: none"> i. No Pockels effect ii. indirect bandgap means that native optical sources are not possible iii. High refractive index means that inherently short devices which are difficult to fabricate (e.g. gratings) iv. Modulation mechanisms tend to be relatively slow v. Thermal effects can be problematic for some optical circuits

Table 2.1: Advantages and disadvantages of silicon photonics over other integrated photonics technologies. Taken from [11].

2.2 Guided-wave photonics

The most important thing for integrated photonics is the way light is confined and manipulated into microscopic structures. While in conventional optics, light is delivered with bulky mirrors and lenses toward the desired position, in integrated photonics waveguides are the main device for this task.

2.2.1 Waveguides

A waveguide is a path inside a medium, whose volume is defined by a certain number of interfaces between different materials, in which light remains confined and ideally travels with negligible losses. This volume is called *core* of the waveguide, while the medium external to it, if present, is called *cladding*.

One can distinguish two types of waveguides: metallic and dielectric. The former are based on the reflection of the electromagnetic field by the metallic surface. Light is confined in the desired path by a series of metallic mirrors. However, this mechanism works well only for electromagnetic radiation for which metals can still be considered *perfect metals*. For visible and infrared frequencies integrated metallic waveguides are not possible, due to too high absorption of the metals at optical frequencies.

On the other hand, dielectric waveguides are based on the phenomenon of total internal reflection (TIR) on the interface between two dielectric media. To produce TIR, the two materials must have sufficiently different real part of the refractive index. Moreover, they must also be transparent, i.e. they must have low imaginary part of the refractive index, in the required range of frequencies. Hence not all media are suitable to build dielectric waveguides that work at a specific wavelength. For example, silicon is known to be transparent for light at and around $1.55\mu\text{m}$, among other wavelengths. This feature together with a low distortion of signals is the reason why it is the most used wavelength in communication technology.

Dielectric Waveguides

Total internal reflection is the well-known phenomenon in which light coming from a material with high refractive index n_H gets reflected at the interface with a material with low refractive index n_L . For this to happen, incident light must be at an angle greater than the critical angle given by the Snell's law:

$$\theta_C = \arcsin\left(\frac{n_L}{n_H}\right).$$

Waveguides have many shapes, however one can initially discriminate them by the dimensionality of the confinement of light. The simplest one is called *slab* and it is composed by two interfaces, which divides the material in three volumes, as shown in Figure 2.2 below. This type is classified as a 1D-waveguide, because it constrains light in only one dimension.

2D-waveguides on the other hand are more common, to the point that with the term waveguide one usually refers to a device belonging to this category. Since they confine light in two dimension, their core then becomes a straight path with which one can conduct light from a place to another. A few types are shown in Figure 2.3 below.

Obviously, the versatility of straight (2D-) waveguides is limited. Hence more complex structures such as bent waveguides have been developed.

Propagation of light inside waveguides

Light propagates inside waveguides as superposition of transverse electro-magnetic (TEM)

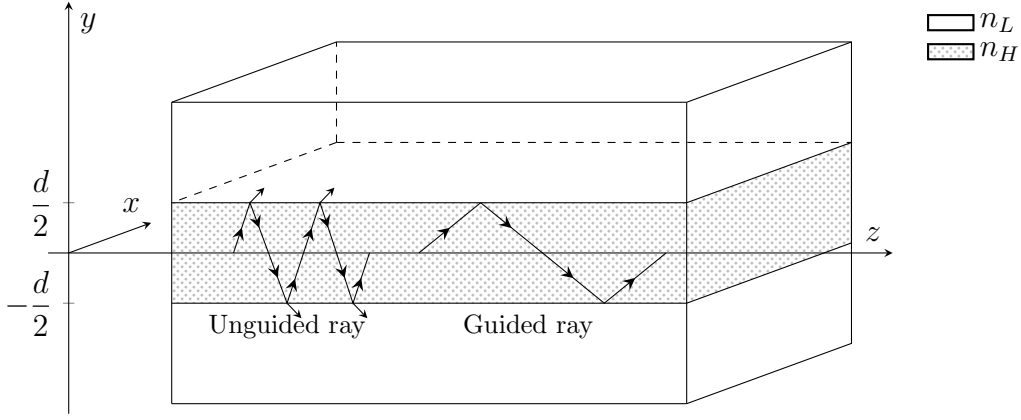


Figure 2.2: Scheme of a dielectric slab waveguide, made by two materials with refractive indexes n_H and n_L , with $n_H > n_L$. Two rays are shown: the unguided one is incident on the interface at an angle smaller than the critical angle. Conversely the guided ray is incident at a greater angle and is therefore totally reflected inside the waveguide [12].

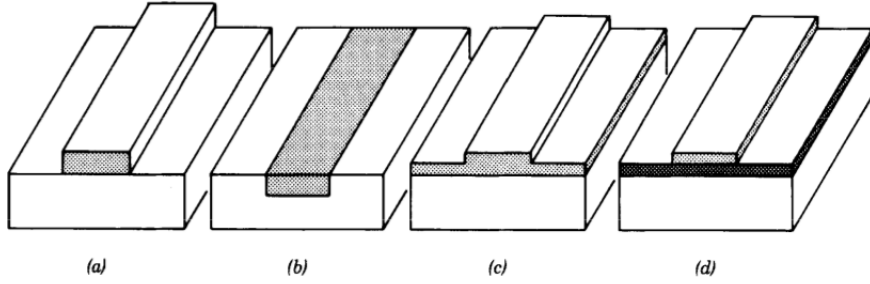


Figure 2.3: Representation of a few types of 2D dielectric waveguides. Light is constrained in two direction and can only move forward or backward in the third direction. (a) strip, (b) embedded strip, (c) rib, (d) strip loaded [12].

waves which keep reflecting on the interfaces of the core. The result of the superposition is the so called *mode*, which is expressed mathematically by

$$\phi(\mathbf{r}) = \phi_m(x, y) e^{i(\beta_m z - \omega t)}. \quad (2.1)$$

$\phi_m(x, y)$ is the transverse field distribution, z is the direction of propagation, β_m is the propagation constant of the m -th mode, and $\omega = 2\pi\nu$ is the angular frequency of light. Each mode, identified by the m index, travels inside the waveguide with a certain propagation constant β and maintaining a field distribution ϕ_m . Both of them depend on the materials and the geometry of the waveguide, however while β_m decrease with increasing m indexes, the features of the function ϕ_m grows in number. Light either propagate in a manner described by a mode or superposition of modes or is rapidly scattered away.

The field distribution of some of the first modes, in the simplest case of a slab waveguide, is represented in Figure 2.4 below. The function inside the core is characterized by maxima, minima and nodes with a certain periodicity and the number of nodes (where the function is zero) is strictly linked to the index m . Outside, the distribution of the field decays exponentially. For waveguide with a 2D core, the field can be expected to be a multivariate version of the same or similar function. However, differently from the analytical solution for slab waveguides, they are almost always found numerically.

The propagation constant can be written as

$$\beta_m = n_{eff} k_0 + i\alpha_{eff}/2 \quad (2.2)$$

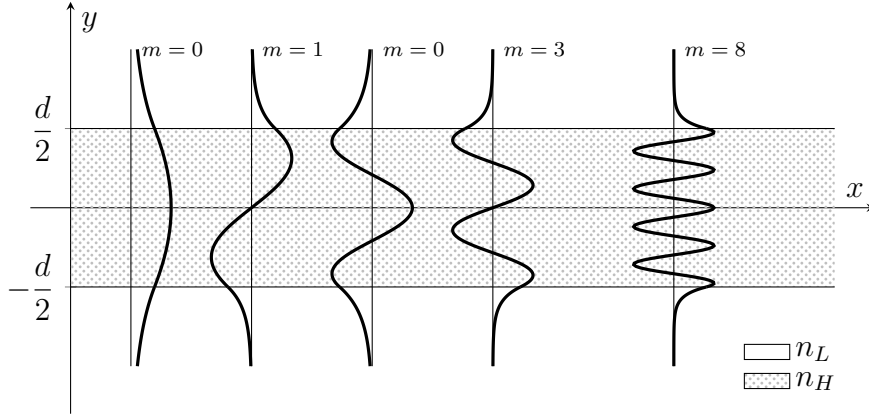


Figure 2.4: Field distribution inside a slab waveguide.

where k_0 is the wavevector in vacuum, $n_{eff} := c_0/v_{ph}$ is the *effective refractive index* of the mode, and α_{eff} is *effective absorption coefficient*. The effective refractive index is the ratio between the speed of light and the *effective phase velocity* v_{ph} at which each wavefront propagates in the core. It is a pure number and its value is between the core refractive index and the cladding refractive index $n_L < n_{eff} < n_H$. A high refractive index means a higher confinement of the electromagnetic field inside the core and a slower propagation speed along the waveguide. The effective absorption coefficient instead is defined by the ratio between the input and output powers for a material of depth L , such that $I_{out}/I_{in} = e^{-\alpha_{eff}L}$, and has units of m^{-1} . A low absorption coefficient is synonym of transparency and low power absorption per unit length.

The physical interpretation of the effective indexes is that they define the propagation of the mode as a whole, instead of a superposition of multiple TEM waves each with a propagation characterized by the local values of refractive index and absorption coefficient. Moreover, in term of magnitudes, a higher refractive index means a slower propagation whereas a higher absorption coefficient means a higher power absorption per unit length.

All of these parameters which describe the propagation of light along the waveguide depends both on the materials and the geometry of the core and the cladding, but also on the frequency of light. For dielectric waveguides, the fundamental mode ($m = 0$) is always supported while higher modes might not be, depending both on the waveguides and on the frequency of light. For this reason one distinguish *single-mode* waveguides, which allow propagation for only the fundamental mode in a certain operative range of frequencies, from *multi-mode* waveguides, which allow higher order modes to propagate.

Equation (2.1) describes the propagation of a monochromatic light wave, which has ideally the same mode amplitude from $t = -\infty$ to $t = +\infty$. This is obviously not a good physical representation, as light is generated and absorbed. Eventually light is described to travel in wavepackets of finite duration, which are inherently non-monochromatic.

When light becomes non-monochromatic, the frequency dependence of the propagation constant β has to be considered. Usually, if working on a restricted range of frequencies, this dependence is characterized by a Taylor expansion at the first order.

$$\beta(\omega) \simeq \beta(\omega') + \frac{\partial \beta(\omega)}{\partial \omega} \Delta \omega \quad (2.3a)$$

$$\simeq \frac{1}{v_{ph}} \omega' + \frac{1}{v_g} \Delta \omega \quad (2.3b)$$

where $1/v_{ph} = n_{eff}/c_0$, $\Delta \omega = \omega - \omega'$, and $v_g := \left(\frac{\partial k}{\partial \omega}\right)^{-1}$. The *group velocity* v_g describes the speed at which the ensemble of frequencies, the wavepacket, propagates. Then, in a similar manner to n_{eff} which is the ratio between c_0 and the effective phase velocity, one can define

the *effective group index* n_{eff}^g as ratio between the speed of light in vacuum c_0 and the effective group velocity v_g of a packet of light inside a waveguide:

$$n_{eff}^g := \frac{c_0}{v_g} = n_{eff}(\omega) + \omega \frac{\partial n_{eff}}{\partial \omega}. \quad (2.4)$$

The latter equation allow to describe the Taylor expansion of β in other terms:

$$\beta(\omega) \simeq \frac{n_{eff}}{c_0} \omega' + \frac{n_{eff}^g}{c_0} \Delta\omega \quad (2.5)$$

Nevertheless the monochromatic wave propagation model is still valid, because it is a very good approximation of slowly varying and almost monochromatic waves, such as the one I will use in my experiments.

GVD?

2.2.2 Microring Optical Cavity

In integrated photonics the microring is an optical cavity made by bending a waveguide on itself. To insert light into the cavity, the microring is often side coupled with one or two straight waveguides. Light inserted from the different ports into the cavity interfere with itself, after a round trip, and with each other. When the waves are in phase, they generate constructive interference and the energy stored in the cavity increases exponentially. On the other hand, when the signals are out of phase, the interference is destructive and the energy stored in the cavity remains low.

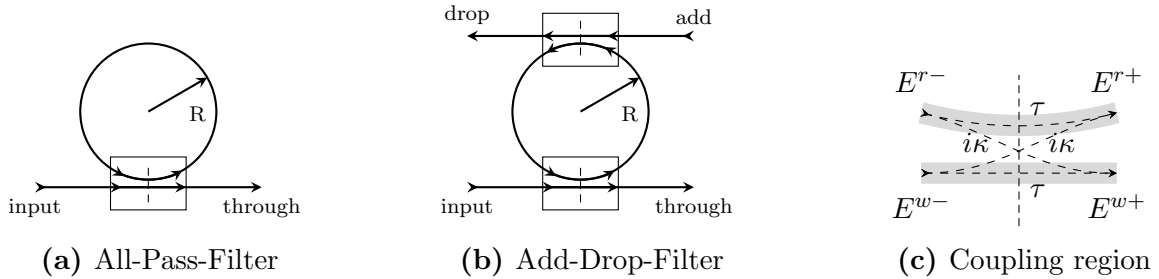


Figure 2.5: Schematic representation of microring configurations and of their coupling regions. APF has two channels, *input* and *through*, and one coupling region. ADF has four channels, *input*, *through*, *add* and *drop*, and two coupling regions.

In the simpler case when there is only one waveguide side coupled to the resonator the optical cavity has only two ports, which are usually called *input* and *through*. This configuration, shown in Figure 2.5a, is called *All-Pass Filter*, because in the ideal case, where no losses happen, all the signal passes from the input to the through channels. All pass resonators have only one input and one output: light then propagate in only in one direction, unless some specific phenomena happens (e.g. back scattering).

By adding a secondary waveguide coupled to the ring, one obtains the so called Add-Drop Filter (ADF) configuration, shown in Figure 2.5b. Its name derive by the fact that the two additional ports are called *add* and *drop* respectively. This simple structure can be readily used as a signal mixer: signals at the resonance wavelengths of the cavity are directed from the input to the drop channel or from the add channel to the through port. Additionally signals out of resonance travel straight from the input to the through channel and from the add to the drop port. Resonators with more than four ports are possible, however are very uncommon.

The theoretical model usually employed to describe resonators analytically decomposes their structure in a series of simpler substructures. Both the APF and the ADF configurations are studied as ordered combinations of straight waveguides and bent waveguides, which are coupled together in specific *coupling regions*, as shown in Figure 2.5c. Other used descriptions, such as numerical simulation obtained with finite element methods (FEM), achieve more precise results, however they are more dependent on the specific geometry of the problem. Moreover the approximation given by the analytical model is sufficiently accurate to make clear the physics behind it and it can be expanded to include non-trivial phenomena (see Subsection 2.2.3). In the following section I will go through the necessary steps to solve the case of the ADF configuration in the theoretical model.

Add-Drop-Filter theory

ADF configuration is obtained when a microring is coupled to two waveguides. Such structure is composed by three different basic structures: four straight waveguides, two bent waveguides which together form the whole ring, and two coupling regions between the waveguides and the microring. Each of these pieces transfers light to or from the outside or another piece.

Since our experiments are carried out in a time scale such that physics phenomena can be considered quasi-static, the field propagating in the device is described by a scalar complex function which, assuming a monochromatic continuous EM wave, loses any temporal dependence:

$$E(z) = |E(z_0)|e^{i\beta z}, \quad (2.6)$$

where z is the direction of propagation, loosely defined to accommodate propagation both along straight and bent waveguides ($z \sim r\theta$), t is time, $\omega = 2\pi\nu$ is the angular frequency, and $\beta = n_{eff}k_0 + i\alpha_{eff}/2$ is the propagation constant. Both dependence of the effective index $n_{eff} = n_{eff}(\omega)$ and the effective loss factor $\alpha_{eff} = \alpha_{eff}(\omega)$ will be expanded only when necessary.

Propagation in the straight waveguides is often considered without loss (i.e. $\alpha_{eff} \approx 0$) and thus neglected. On the other hand, propagation in the two halves of the ring is frequently considered with radiative losses (i.e. $\alpha_{eff} > 0$), because a bent waveguide is intrinsically more difficult to fabricate in comparison to a straight one. The last structure type, the coupler between the resonator and the waveguides, is considered alike a beamsplitter. In this approximation, light does not propagate in this part as it was in the previous ones. Its operation is reduced to the exchange of power between the two input ports and two output ports, which is described by the following matrices:

$$\begin{pmatrix} E_{ch}^{w+} \\ E_{ch}^{r+} \end{pmatrix} = \mathbf{M} \begin{pmatrix} E_{ch}^{w-} \\ E_{ch}^{r-} \end{pmatrix}, \quad \mathbf{M} = \begin{pmatrix} \tau & i\kappa \\ i\kappa & \tau \end{pmatrix}, \quad (2.7)$$

where E is the field amplitude of the ‘ $_{ch}$ ’ channel either in the waveguide ‘ w ’ or in the ring ‘ r ’, before ‘ $-$ ’ and after ‘ $+$ ’ the coupling. An explanatory diagram is shown in Figure 2.5c. The matrix \mathbf{M} is characterized by two real valued parameters, τ and κ , between 0 and 1. Specifically, they must verify the following constraint:

$$\det(\mathbf{M}) = |\tau|^2 + |i\kappa|^2 = 1. \quad (2.8)$$

which represents the conservation of energy.

At this point one can solve the problem by putting all the pieces together. Neglecting the

propagation inside the waveguide, the full system of equation that describe the problem is:

$$E_{in}^{w+} = \tau E_{in}^{w-} + i\kappa E_{in}^{r-} \quad (2.9a)$$

$$E_{in}^{r+} = i\kappa E_{in}^{w-} + \tau E_{in}^{r-} \quad (2.9b)$$

$$E_{out}^{r-} = E_{in}^{r+} e^{i\beta\pi R} \quad (2.9c)$$

$$E_{in}^{r-} = E_{out}^{r+} e^{i\beta\pi R} \quad (2.9d)$$

$$E_{out}^{w+} = \tau E_{out}^{w-} + i\kappa E_{out}^{r-} \quad (2.9e)$$

$$E_{out}^{r+} = i\kappa E_{out}^{w-} + \tau E_{out}^{r-} \quad (2.9f)$$

The first two equations (2.9a) and (2.9b) describe the exchange of optical power between the first channel and the microring resonator. Then the third and fourth equations (2.9c) and (2.9d) delineate the propagation of light in the two halves of the resonator, from one coupling region to the other. Lastly the remaining two equation characterize the transfer of light between the resonator and the add-drop channel.

Using these equations it is easy to define new quantities of interest: the first one is transmittance from the *input* to the *through* port.

$$\eta_T(\omega) := \frac{E_{in}^{w+}}{E_{in}^{w-}} = t \frac{1 - e^{i\beta 2\pi R}}{1 - \tau^2 e^{i\beta 2\pi R}} \quad (2.10)$$

Similarly, one can also define the transmittance from the *input* to the *drop* port.

$$\eta_D(\omega) := \frac{E_{out}^{w+}}{E_{in}^{w-}} = \frac{-\kappa^2 e^{i\beta\pi R}}{1 - \tau^2 e^{i\beta 2\pi R}} \quad (2.11)$$

However, this quantities have complex values and therefore are difficult to study. For this reason, one define also the transmission between the same ports as the square modulus of the transmittance. Thus it follows that

$$T(\omega) := |\eta_T|^2 = \tau \frac{(1 - \gamma)^2 + 4\gamma \sin^2(n_{eff} k_0 \pi R)}{(1 - \tau^2 \gamma)^2 + 4\tau^2 \gamma \sin^2(n_{eff} k_0 \pi R)} \quad (2.12)$$

and

$$D(\omega) := |\eta_D|^2 = \frac{\kappa^4 \gamma}{(1 - \tau^2 \gamma)^2 + 4\tau^2 \gamma \sin^2(n_{eff} k_0 \pi R)} \quad (2.13)$$

where the notation is simplified by loss parameter $\gamma = e^{-\alpha_{eff} \pi R}$. The dependence of $T = T(\omega)$ and $D = D(\omega)$ from the frequency (or wavelength) of light is obtained by making explicit the wavevector dependence $k_0 = \frac{\omega}{c_0} = \frac{2\pi}{\lambda}$ in which c_0 is the universal constant speed of light in vacuum.

The transmission is much more interesting to study than the transmittance because it represents the ratio between the input and output optical powers ($I \propto |E|^2$), except for a constant factor, and is therefore a real valued function of ω instead of a complex one. Both transmission spectra, as shown in Figure 2.6, have either peaks or dips: the resonances of the microring. Specifically, at each resonance, the through spectrum shows dips while the drop spectrum shows peaks. The depth of the dips, the height of the peaks, and the width of both of them is defined by few parameters: the coupling constant τ^1 and the half round trip loss factor $\gamma = e^{-\alpha_{eff} \pi R}$.

¹one can choose also κ , but it does not matter which one, since $\tau^2 = 1 - \kappa^2$

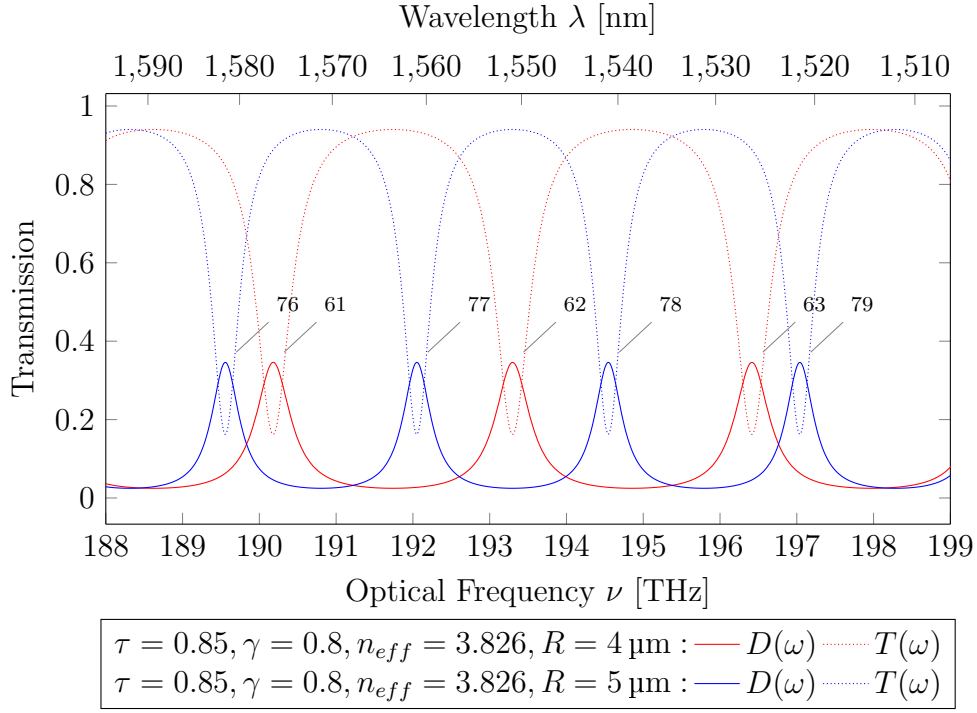


Figure 2.6: Transmission spectra of microring resonators in Add-Drop Filter configuration for the through and drop ports. The resonators considered differ in radius ($R=4\mu\text{m}$ and $R=5\mu\text{m}$), but share the coupling constant $\tau = 0.85$, the half round trip loss factor $\gamma = 0.8$ and the effective index $n_{eff} = 3.826$.

The frequency of each resonance is identified by a positive integer number, that verifies the following equation:

$$\omega_m = m \frac{c_0}{n_{eff}(\omega)R}, \quad (2.14)$$

where the dependence of $n_{eff} = n_{eff}(\omega)$ has been explicitly shown. Moreover, it is easy to convert this quantity in other domains with:

$$\omega = 2\pi\nu \quad \text{and} \quad c_0 = \nu\lambda \quad (2.15)$$

Each resonance is spaced from the next one by a quantity called *Free Spectral Range* FSR_{ω_m} . By exploiting the Taylor expansion of β seen in equation (2.5) one obtains

$$FSR_{\omega_m} \simeq \frac{c_0}{n_{eff}^g(\omega_m)R}. \quad (2.16)$$

where the *effective group index* $n_{eff}^g(\omega)$ has been defined in equation (2.4). Another important quantity is the width of the peaks or dips of each resonance. Ordinarily one evaluate the so called *full-width-half-maximum* (FWHM):

$$FWHM_{\omega_m} = \frac{c_0}{n_{eff}^g(\omega_m)} \frac{1 - \tau^2\gamma}{\pi R \tau \sqrt{\gamma}} \quad (2.17)$$

Values such this are often expressed in the wavelength domain: however to obtain the FSR_λ one can not use equation (2.15). Instead the following relations have to be used:

$$\Delta\omega = 2\pi\Delta\nu = \frac{2\pi c_0}{\lambda^2}\Delta\lambda \quad \text{and} \quad \Delta\lambda = \frac{c_0}{\nu^2}\Delta\nu = \frac{2\pi c_0}{\omega^2}\Delta\omega \quad (2.18)$$

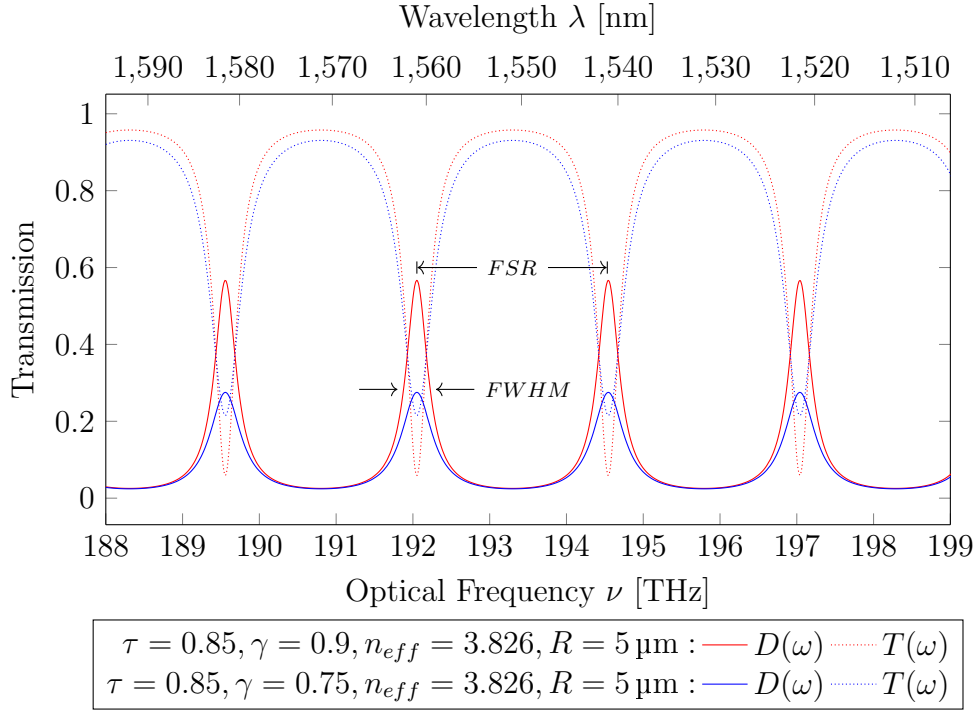
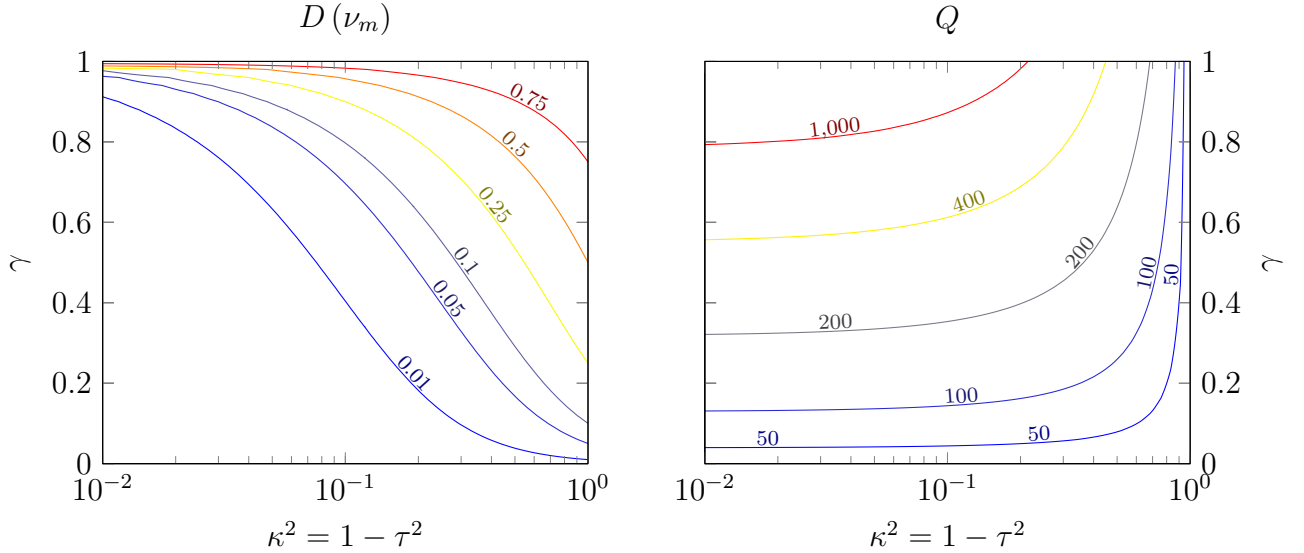


Figure 2.7: Transmission spectra of microring resonators in Add-Drop Filter configuration for the through and drop ports. The microring has a radius of $5 \mu\text{m}$ and an effective index $n_{eff} = 3.826$. The coupling constant is $\tau = 0.85$, while the half round trip loss factor γ takes the two values of 0.9 and 0.75. The arrows indicate a $FSR_\nu \approx 2.49 \text{ THz}$ and a $FWHM_\nu \approx 0.35 \text{ THz}$, which gives a quality factor of around $Q \approx 550$. The other resonator has $FWHM_\nu \approx 0.5 \text{ THz}$ which gives a quality factor of $Q \approx 390$.



By employing the first relation, one obtains

$$FSR_{\lambda_m} \simeq \frac{\lambda_m^2}{n_{eff}^g(\lambda_m) 2\pi R} \quad \text{and} \quad FWHM_{\lambda_m} \simeq \frac{\lambda_m^2}{n_{eff}^g} \frac{1 - \tau^2 \gamma}{2\pi^2 R \tau \sqrt{\gamma}} \quad (2.19)$$

where $n_{eff}^g(\lambda_m)$ is an appropriate redefinition of the effective group index in the wavelength domain.

An additional figure of merit of an optical microring resonator is the *quality factor* Q – *factor*, or simply Q . The quality factor has many similar definitions, depending on the field of study. In the most physical, but less operative one defines it as 2π times the ratio between the energy stored in the cavity and the energy lost each cycle.

$$Q := 2\pi \times \frac{\text{energy stored}}{\text{energy lost per cycle}}$$

However, for my purposes this definition is too cumbersome, thus a similar one is used instead:

$$Q := \frac{\omega_m}{FWHM_{\omega_m}} = \frac{\lambda_m}{FWHM_{\lambda_m}} \quad (2.20)$$

Both are sensible definitions and as Q becomes larger they become approximately equivalent.

Critical coupling

Critical coupling/ over-coupling / under-coupling / enhancement factor

We covered how light travels inside a waveguide, now the problem on how to insert and extract light from it. The process of inserting light into a waveguide is called *coupling* and it often requires at the same time to extract light from another waveguide.

The most naive way to insert light into a waveguide is obviously to use one of its ends. This is achieved either by radiating a beam directly on the whole waveguide or by focusing it inside the core with some lensing mechanism (e.g. tapered optical fibers).

These methods are very simple, however they have a low coupling efficiency.

A more sophisticated way to couple light inside a waveguide is to use the so called *grating coupler*. A grating coupler is a periodic structure end coupled to a waveguide. It is designed in a way such that light incident from free space above it, at a certain angle and with a specific wavelength, is coupled inside the waveguide.

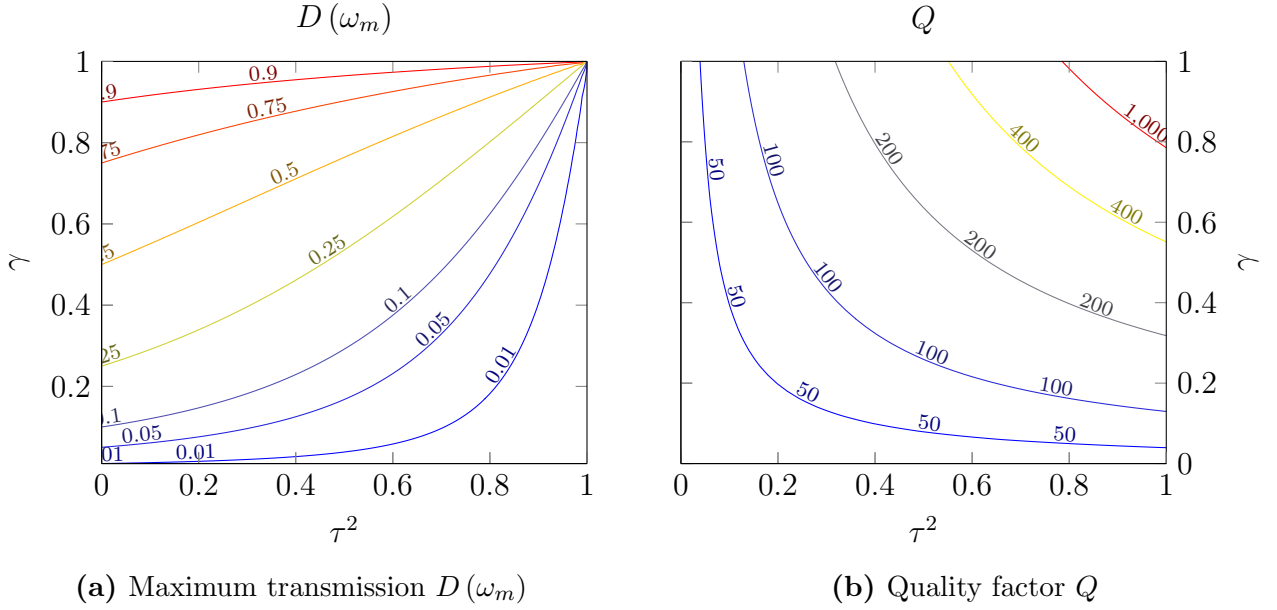


Figure 2.9: Contour plots of the maximum transmission on the drop channel $D(\omega_m)$ and of the quality factor Q , parametrized for the coupling coefficient τ^2 and the half round trip loss factor γ . The Q – factor is obtained for a resonance near $\nu_m \approx 200$ THz (arbitrarily chosen) and an effective group index $n_{eff}^g \approx 4$.

Evanescent coupling exploits the fact that the evanescent field distribution outside the core decays exponentially to zero [11]. This means that when two waveguides are separated by a sufficiently small distance, the field distribution of one waveguide cannot decay sufficiently fast to zero and instead extends itself over the core of the second waveguide. When this happens optical power can be transferred between the waveguides. Therefore, if at the beginning light is propagating inside one waveguide, light is *coupled* inside the other waveguide because its mode overlaps the core of the second waveguide. Depending on the parameters of the problems, light can couple completely or partially between the waveguides.

This exchange of optical power is periodic with their length.

2.2.3 Nonlinear Perturbations

In the precedent section, I considered the material of waveguides and resonators as a linear medium. This is however only an approximation, because almost any material shows some sort of nonlinearity if probed with a high enough electromagnetic field.

In integrated photonics, due to the lateral confinement of the field inside the waveguide, much higher field amplitude are reached in comparison to free space. Moreover, inside an optical cavity, such as an appropriately built microring resonator, the field reaches even more higher amplitudes thanks to its large enhancement factor and nonlinearities might arise.

Silicon response to light is almost linear. However, a silicon optical cavity, by confinement and enhancement, can obtain a very high electromagnetic field inside, such that its response becomes significantly nonlinear.

One can distinguish two kinds of optical nonlinearities shown by silicon: electronic nonlinearities and thermal nonlinearities.

Electronic nonlinearities

Silicon is a centrosymmetric material and thus does not exhibit optical nonlinearities of even orders. The most efficient nonlinear effects belong to the third order.

Silicon nonlinearities are mainly due to two fundamental effects of this category and their correlated effects.

Kerr effect The first one is the Kerr effect, which is given by the real part of the third order nonlinear susceptibility $\chi^{(3)}$ and it produces a change in the refractive index characterized by:

$$\Delta n_{Kerr} = n_2 I, \quad (2.21)$$

where I is the intensity of light beam and n_2 is the second-order (or intensity dependent) nonlinear refractive index, defined by

$$n_2 = \frac{3}{4\varepsilon_0 n_0^2 c_0} \text{Re} [\chi^{(3)}] \stackrel{Si}{\simeq} 0.45 \times 10^{-17} \text{ m}^2 \text{ W}^{-1}. \quad (2.22)$$

Two photon absorption The other fundamental effect is the two photon absorption (TPA), which is linked instead to the imaginary part of the third order nonlinear susceptibility. It does not produce a change to the real part of the refractive index but to its imaginary part, which is in turn linked to the absorption coefficient by $\alpha = 2\text{Im}[n]\omega/c_0$:

$$\Delta \alpha_{TPA} = \beta_{TPA} I, \quad (2.23)$$

where β_{TPA} is the TPA coefficient which is experimentally found to be

$$\beta_{TPA} \stackrel{Si}{\simeq} 0.79 \times 10^{-11} \text{ m W}^{-1}. \quad (2.24)$$

As a consequence to the absorption of two photons, a free carrier is generated in the conduction band of silicon. The presence of free carriers alters the refractive index in its real and imaginary parts as well with two additional effects.

Free carrier effects When free carriers are generated in the conduction band, they can either move around until they thermalize or they can absorb an incoming photon. The first case is the *free carrier dispersion* (FCD) effect, which affects the real part of the refractive index. The second case is the *free carrier absorption* (FCA) effect and affects the imaginary part of the refractive index. Both effects are usually characterized as first order expansion, as follows:

$$n(N) = n(N_0) + \left. \frac{dn}{dN} \right|_{N_0} \Delta N \quad (2.25)$$

$$\alpha(N) = \alpha(N_0) + \left. \frac{d\alpha}{dN} \right|_{N_0} \Delta N \quad (2.26)$$

where the value of the FCD and FCA coefficients is

$$\left. \frac{dn}{dN} \right|_{N_0} \stackrel{Si}{\simeq} -1.73 \times 10^{-21} \text{ m}^3 \quad \text{and} \quad \left. \frac{d\alpha}{dN} \right|_{N_0} \stackrel{Si}{\simeq} 1.1 \times 10^{-15} \text{ m}^2 \quad (2.27)$$

respectively.

Thermal nonlinearities

The thermal nonlinearity of materials refractive index linked to their optical response is called *thermo-optic effect* (TOE). When light propagates inside a medium, a portion of the photons is absorbed. Hence the material heats up and its refractive index changes. To characterize this change usually a first order expansion is made:

$$n(T) = n(T_0) + \left. \frac{dn}{dT} \right|_{T_0} \Delta T, \quad (2.28)$$

where $\frac{dn}{dT}$ is called *thermo-optic coefficient* and ΔT is the temperature shift.

In silicon waveguides, the shift in temperature is caused both by linear and nonlinear processes. The linear process is characterized by the linear absorption coefficient α . Similarly, the nonlinear processes of heat generation are linked to the nonlinear part of the absorption coefficient and are given mainly by the two photon absorption and by the free carrier absorption.

The thermo-optic coefficient (TOC) of silicon at 300 K is ?? :

$$\left. \frac{dn}{dT} \right|_{300\text{ K}} = 1.86 \times 10^{-4} \text{ K}^{-1}. \quad (2.29)$$

Moreover, due to the fact that in silicon integrated structures the cladding is usually made by SiO_2 , a thermally insulating material, the heat generated by the beam confined in the cavity will stay within the core, thus amplifying the effect even further.

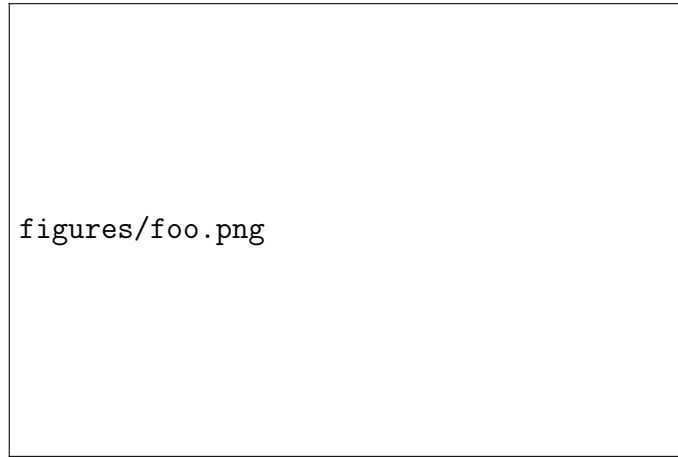


Figure 2.10: •

2.3 Integrated photonics applied to ANNs

Since our experiments are carried out in a time scale such that physics phenomena can be considered quasi-static, we obtain also that the only nonlinear effect non negligible is the thermo-optic effect.

2.3.1 Weighted sum of inputs

This has already been demonstrated and integrated widely, so it will not be the focus of this work. Two example are the banks of microring resonators and the MZ interferometers.

2.3.2 Nonlinear Activation Function

As opposed to the mechanism for weighted sum, an optical phenomenon for the activation function in an integrated photonic circuit has yet to be proposed.

2.3.3 Simulations

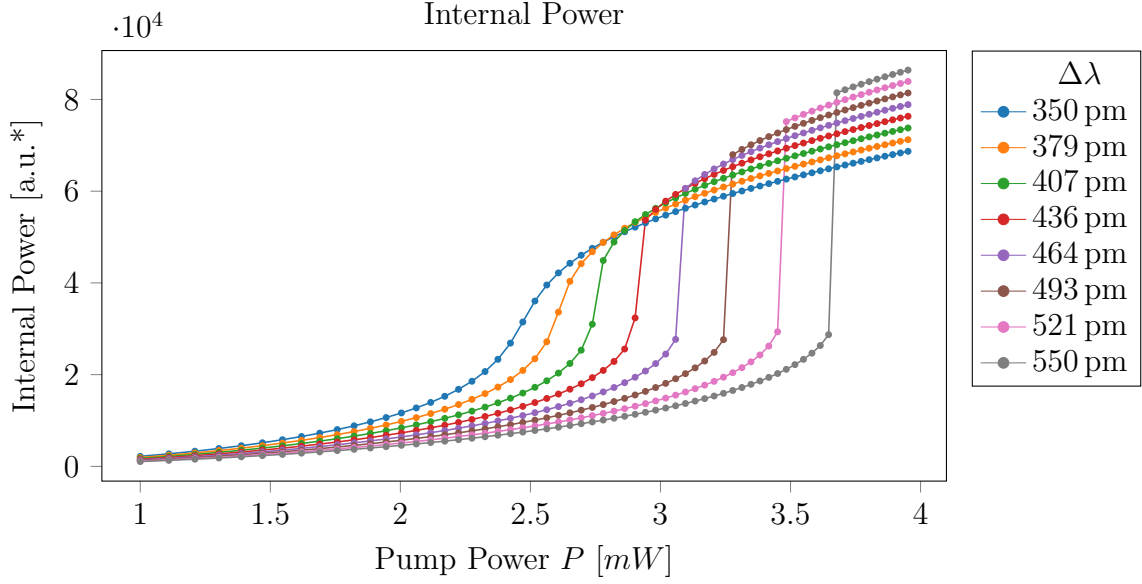


Figure 2.11

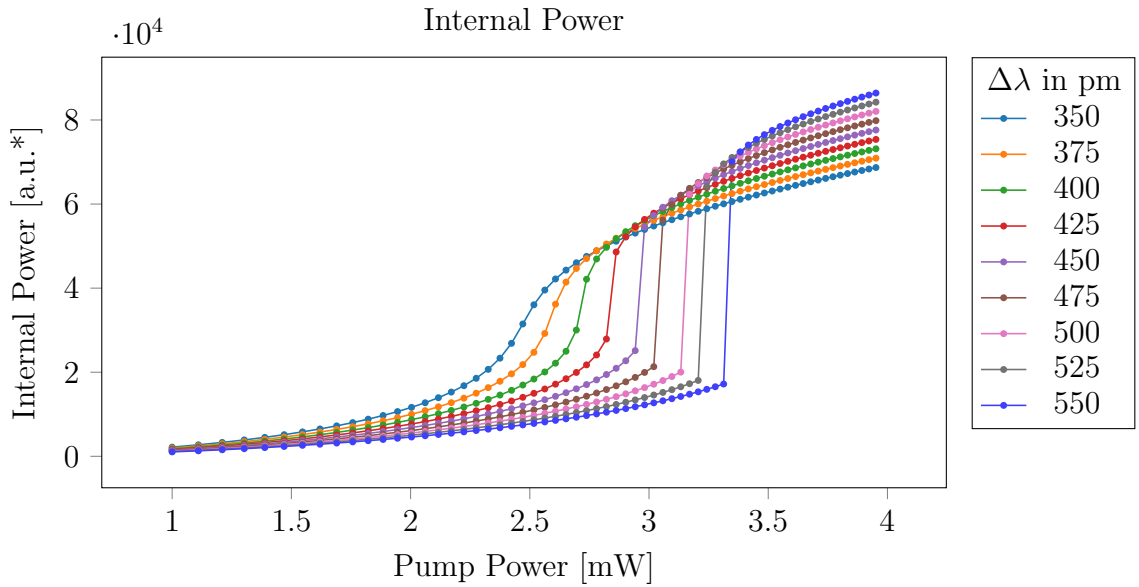


Figure 2.12

Chapter 3

Samples, setup and experiments

All the experiments have been made on samples manufactured for the IRIS project. This is due to the fact that in the time frame of this work there would have not been enough time to design and produce an ad hoc device. Moreover, as already stated, the aim of this thesis is to produce a proof of concept for an all-optical implementation of an activation function, rather than to construct a complete prototype.

3.1 The samples

The IRIS project studied the design and the implementation of an integrated reconfigurable silicon photonics switch matrix, a routing device, as a replacement for electronic devices used in the telecommunication industry. The completed integrated photonic circuit consists of a matrix of waveguides crossing each other and linked by couples of racetrack resonators, thermally-controlled. At the ends of the waveguides other structures, interleavers and AWGs, allowed many signals at different wavelengths to be multiplexed/demultiplexed on/from the same waveguide.

The complexity of such photonic circuit required the fabrication, for testing purposes, of each and every structure of which it is composed, with various production parameters. The collection of all these devices on a chip was produced in many samples. Isolated from the others, but also grouped together with few other elements. Specifically these test structures were disposed on a single chip and were accessible via grating couplers.

Since this work is at the beginning of the project, my choice was a system of intermediate complexity, so that it could be used both to test the activation function and the weighted sum.

The structure selected is the following: a simple waveguide, coupled to eight drop channels by a single or a couple of ring resonators each, nicknamed *mini-matrix*. In this family of devices, there were those built with single microrings, double microrings, single racetracks, or double racetracks. The final choice was to study the *mini-matrix* in which the coupling mechanism was provided by single ring resonators, because of its simpler transfer function in respect to double microrings or double racetracks.

Moreover, these samples were manufactured with thermo-electric **this is not the correct term** pads to heat the rings and effectively tune their resonance.

3.2 Setup

3.3 Characterization of the Activation Function

Characterization of the activation function

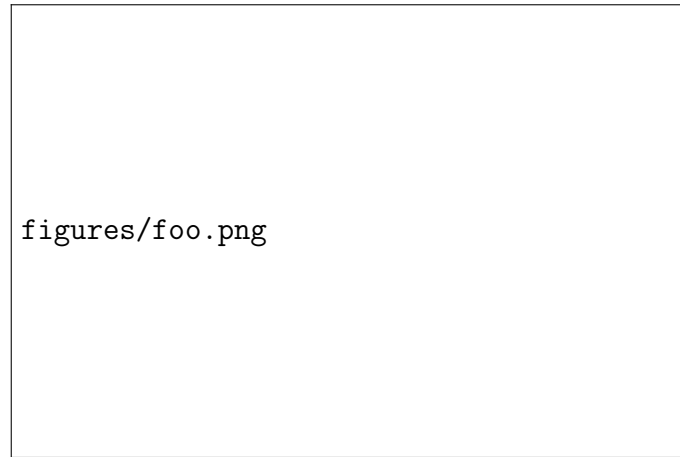


Figure 3.1: image/scheme of the minimatrix

3.4 Test of a Trained ANN

Test of the activation function

Conclusion

Furthermore, regarding the second objective, my work is part of a bigger project, BACUKP, with a wider scope and a more ambitious goal. It attempts to connect three fields of research, normally quite distant from each other: physics, computing technology, and biology. The main idea is to use the integrated photonics framework (physics) to build an artificial neural network (computing technology) on chip, where one could grow and, at least partially control, real neurons (biology). The aim is to develop a methodologies and devices that will allow us to study neural activity from a bottom-up perspective.

Bibliography

- [1] *Humanbrainproject.eu*. [Online]. Available: <https://www.humanbrainproject.eu/en/silicon-brains/>.
- [2] W. S. McCulloch and W. Pitts, “A logical calculus of the ideas immanent in nervous activity”, *The Bulletin of Mathematical Biophysics*, vol. 5, no. 4, pp. 115–133, 1943, ISSN: 00074985. DOI: [10.1007/BF02478259](https://doi.org/10.1007/BF02478259). arXiv: [arXiv:1011.1669v3](https://arxiv.org/abs/1011.1669v3).
- [3] D. O. Hebb *et al.*, *The organization of behavior: A neuropsychological theory*, 1949.
- [4] B. Farley and W. Clark, “Simulation of self-organizing systems by digital computer”, *IRE Professional Group on Information Theory*, vol. 4, no. 4, pp. 76–84, 1954, ISSN: 2168-2690. DOI: [10.1109/TIT.1954.1057468](https://doi.org/10.1109/TIT.1954.1057468). [Online]. Available: <http://ieeexplore.ieee.org/lpdocs/epic03/wrapper.htm?arnumber=1057468>.
- [5] N. Rochester, J. H. Holland, L. H. Haibt, and W. L. Duda, “Tests on a cell assembly theory of the action of the brain, using a large digital computer”, *IRE Transactions on Information Theory*, vol. 2, no. 3, pp. 80–93, 1956, ISSN: 21682712. DOI: [10.1109/TIT.1956.1056810](https://doi.org/10.1109/TIT.1956.1056810).
- [6] F. Rosenblatt, “The perceptron a perceiving and recognizing automaton”, *tech. rep., Technical Report 85-460-1*, 1957.
- [7] F. Rosenblatt, “The perceptron: A probabilistic model for information storage and organization in the brain”, *Psychological Review*, vol. 65, no. 6, pp. 386–408, 1958, ISSN: 0033295X. DOI: [10.1037/h0042519](https://doi.org/10.1037/h0042519). arXiv: [arXiv:1112.6209](https://arxiv.org/abs/1112.6209). [Online]. Available: <http://psycnet.apa.org/journals/rev/65/6/386.pdf%7B%5C%%7D5Cnpapers://c53d1644-cd41-40df-912d-ee195b4a4c2b/Paper/p15420>.
- [8] M. Minsky and S. Papert, “Perceptrons: An Introduction to Computational Geometry”, *MIT Press*, p. 268, 1969, ISSN: 00189340. DOI: [10.1109/T-C.1969.222718](https://doi.org/10.1109/T-C.1969.222718).
- [9] C. Clabaugh, D. Myszewski, and J. Pang, *Neural networks - history*, 2000. [Online]. Available: <https://cs.stanford.edu/people/eroberts/courses/soco/projects/neural-networks/History/index.html>.
- [10] *Pytorch.org*. [Online]. Available: <http://pytorch.org/about/>.
- [11] G. T. Reed and A. P. Knights, *Silicon Photonics*. 2004, pp. 166–170, ISBN: 9780470014189. DOI: [10.1002/0470014180](https://doi.org/10.1002/0470014180). [Online]. Available: <http://doi.wiley.com/10.1002/0470014180>.
- [12] B. E. a. Saleh, M. C. Teich, S. Editor, and J. W. Goodman, *Fundamentals of Photonics (Wiley Series in Pure and Applied Optics)*. 1991, vol. 5, p. 992, ISBN: 0471839655. DOI: [10.1002/0471213748.ch1](https://doi.org/10.1002/0471213748.ch1).

Acknowledgements