



# Bioinformatics LAB 4

## Gene Fusions



**Prof.ssa Elisa Ficarra**

**Prof.ssa Santa Di Cataldo**

**Eng. Marta Lovino**

**Eng. Alessio Mascolini**

Politecnico di Torino

DAUIN

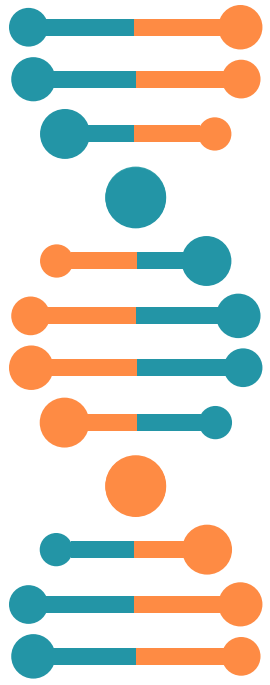
Dept. of Control and Computer Engineering



# LAB 4 - Objectives

# Objectives

- Understand how gene fusions work
- Build the gene fusion sequence starting from its breakpoints
- (Advanced, for the next labs): mono-dimensional CNN for gene fusions prioritization





The background features a light blue DNA double helix on the left side. Scattered across the teal and white background are several chemical structures: a large polycyclic aromatic hydrocarbon (PAH) in the center-left, a benzene ring with a substituent at the bottom right, and several smaller molecules like water (H<sub>2</sub>O) and carbon dioxide (CO<sub>2</sub>) in the upper right. A thick teal diagonal band runs from the top right towards the bottom left.

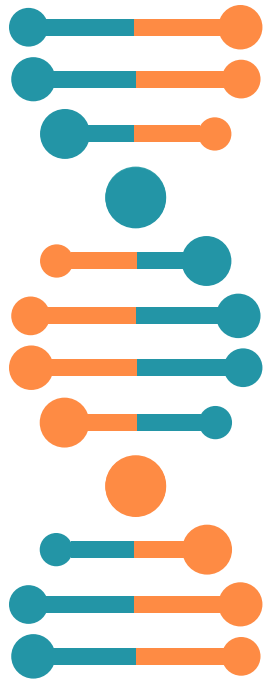
# LAB 3 - Assignments

# Assignment 1: Build the gene fusion sequences

Build the gene fusion sequences starting from the breakpoints.txt file and the reference genome on chromosomes 10 and 18:

breakpoints.txt

Gene5p	chr5p	breakpoint5p	strand5p	Gene3p	chr3p	breakpoint3p	strand3p
ABL1	10	15566488	+	PACX3	18	8666148	+
PPCHS	18	747292	+	TTTCS	10	7393	-
CSSP	10	9999845	-	PPAJD	10	6628	-



# Assignment 2: Perform gene fusion prioritization using CNN

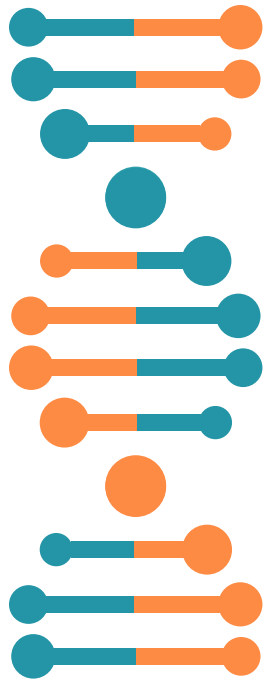
Use the training\_set.csv and test\_set\_1.csv files to train and test a 1-D CNN model to recognize oncogenic and not oncogenic gene fusions. The files are uploaded on the Teaching Portal. Please refer to the recording for all the details.

	FusionPair Chr3p	Label Coord3p	Version 3pStrand	Chr5p 3pCommonName	Coord5p 3pEnsg	5pStrand 3pGeneFunctionality	5pCommonName 3pGeneDescription	5pEnsg MainProteins	5pGeneFunctionality Proteins	5pGeneDescription
0	CSNK2B_NDUFA6 polypeptide [Source:HGNC Symbol;Acc:2460] [Source:HGNC Symbol;Acc:7690]	1 22	grch37 42486683	6 -	31637695 NDUFA6	+ ENSG00000184983	CSNK2B protein_coding	ENSG00000204435 NADH dehydrogenase (ubiquinone) 1 alpha subcomplex, 6, 14kDa	protein_coding	casein kinase 2, beta
	['MEKCKGTSRMAGTTSADVKMSSSEEVSWISWFCGLRGNEFFCEVDEDIQDKFNLTGLNEQVPHYRQALDMILDLEPDEELEDNPNQSDLIEQAAEMLYGLIHARYILTNRGIAQMLEKYQQGDFGYCPRVYCNQPMPLIGLSDIPGEAMVKLYCPKCMDVYTPKSSRHHHTDGAYFGTGFPHMLF MVHPEYRPKRPANQFVPRLYGFKIHPMAYQLQLQAASNFSPVKTIGT'] ['MEKCKGTSRMAGTTSADVKMSSSEEVSWISWFCGLRGNEFFCEVDEDIQDKFNLTGLNEQVPHYRQALDMILDLEPDEELEDNPNQSDLIEQAAEMLYGLIHARYILTNRGIAQMLEKYQQGDFGYCPRVYCNQPMPLIGLSDIPGEAMVKLYCPKCMDVYTPKSSRHHHTDGAYFGTGFPHMLF MVHPEYRPKRPANQFVPRLYGFKIHPMAYQLQLQAASNFSPVKTIGT', 'MEKCKGTSRMAGTTSADVKMSSSEEVSWISWFCGLRGNEFFCEVDEDIQDKFNLTGLNEQVPHYRQALDMILDLEPDEELEDNPNQSDLIEQAAEMLYGLIHARYILTNRGIAQMLEKYQQGDFGYCPRVYCNQPMPLIGLSDIPGEAMVKLYCPKCMDVYTPKSSRHHHTDGAYFGTGFPHMLF YGFKIHPMAYQLQLQAASNFSPVKTNGTG', 'MSSSEEVSWISWFCGLRGNEFFCEVDEDIQDKFNLTGLNEQVPHYRQALDMILDLEPDEELEDNPNQSDLIEQAAEMLYGLIHARYILTNRGIAQMLEKYQQGDFGYCPRVYCNQPMPLIGLSDIPGEAMVKLYCPKCMDVYTPKSSRHHHTDGAYFGTGFPHMLF FKSPVKTIGT', 'MSSSEEVSWISWFCGLRGNEFFCEVDEDIQDKFNLTGLNEQVPHYRQALDMILDLEPDEELEDNPNQSDLIEQAAEMLYGLIHARYILTNRGIAQMLEKYQQGDFGYCPRVYCNQPMPLIGLSDIPGEAMVKLYCPKCMDVYTPKSSRHHHTDGAYFGTGFPHMLF FKSPVKTNGTG']									
1	IFT57_CALM1 57 homolog (Chlamydomonas) [Source:HGNC Symbol;Acc:17367] [Source:HGNC Symbol;Acc:1442]	1 14	grch37 14	3 90871138	107884314 +	- CALM1	IFT57 ENSG00000198668	ENSG00000114446 protein_coding	protein_coding	intraflagellar transport calmodulin 1 (phosphorylase kinase, delta)
	['MTAALAVVTTSGLEDGVPFRSGEGTGEVVLERGPGAAYHMFVVMEDLVEKLKLLRYEEFLRKS NLKAPSRHYFALPTNPGEQFYMFTCTLAAWLINKAGRPFEPQEQYDDPNATISNLSRLSFGRTADFPSPKLKSGYGEHVCYVLDCAEEALYIGFTWKRIPIYPVEEEESVAEDDAELTNKLVDE EFVEEETDNEENFIDLNLVKAQTYHLDNMNETAKQEDILESTTDAAEWSLEVERVLPLQKVITRTDNKDWRIHVDQMHHQHRSGIESALKETKGFGLDKLHNEITRTLEKISSREKYINNQLNLVQYERAAQAQLSEAKERYQQGNL'] ['MTAALAVVTTSGLEDGVPFRSGEGTGEVVLERGPGAAYHMFVVMEDLVEKLKLLRYEEFLRKS NLKAPSRHYFALPTNPGEQFYMFTCTLAAWLINKAGRPFEPQEQYDDPNATISNLSRLSFGRTADFPSPKLKSGYGEHVCYVLDCAEEALYIGFTWKRIPIYPVEEEESVAEDDAELTNKLVDE EFVEEETDNEENFIDLNLVKAQTYHLDNMNETAKQEDILESTTDAAEWSLEVERVLPLQKVITRTDNKDWRIHVDQMHHQHRSGIESALKETKGFGLDKLHNEITRTLEKISSREKYINNQLNLVQYERAAQAQLSEAKERYQQGNL', 'MTAALAVVTTSGLEDGVPFRSGEGTGEVVLERGPGAAYHMFVVMEDLVEKLKLLRYEEFLRKS NLKAPSRHYFALPTNPGEQFYMFTCTLAAWLINKAGRPFEPQEQYDDPNATISNLSRLSFGRTADFPSPKLKSGYGEHVCYVLDCAEEALYIGFTWKRIPIYPVEEEESVAEDDAELTNKLVDEEFVEEETDNEENFIDLNLVKA QTYHLDNMNETAKQEDILESTTDAAEWSLEVERVLPLQKVITRTDNKDWRIHVDQMHHQHRSGIESALKETKGFGLDKLHNEITRTLEKISSREKYINNQLNLVQYERAAQAQLSEAKERYQQGNL', 'MTAALAVVTTSGLEDGVPFRSGEGTGEVVLERGPGAAYHMFVVMEDLVEKLKLLRYEEFLRKS NLKAPSRWVPASPARVPASPQGDALLGPAGNQFRESGHYFALPTNPGEQFYMFTCTLAAWLINKAGRPFEPQEQYDDPNATISNLSRLSFGRTADFPSPKLKSGYGEHVCYVLDCAEEALYIGFTWKRIPIYPVEEEESVAEDDAELTNKLVDEEFVEEETDNEENFIDLNLVKA QTYHLDNMNETAKQEDILESTTDAAEWSLEVERVLPLQKVITRTDNKDWRIHVDQMHHQHRSGIESALKETKGFGLDKLHNEITRTLEKISSREKYINNQLNLVQYERAAQAQLSEAKERYQQGNL']									



# LAB4 – Take home message

- The most common configuration for a gene fusion is “promoter-end”, although multiple configurations are possible (refer to the recording).
- Classify gene fusions is usually really challenging: an algorithm that works well on a training set can fail on the test set.





Questions?

Remember:  
no question is  
stupid