# A Semantic Web approach to data integration for the histone code case

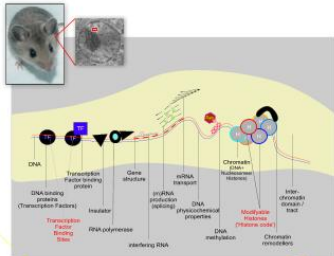M. Scott Marshall

Marco Roos

Lennart JG Post

Timo M. Breit

*Acknowledgements*

Roel van Driel

Paul Fransz

Willem van Hage

Machiel Jansen

Our objectives

1. to enable data integration experiments in terms of our biological knowledge with Semantic Web formats and tools

2. to provide a basis for investigating semantic models as a tool for understanding biological systems

**'Computational experiment'**

```python
# (c) 2006, SARA, Bart Heupers
# Example use :
# python get_overlap.py TFBSConsSites_chr12.txt encode_regions_chr12.txt
> overlaps.txt
import sys, os

def read_line(f):
    '''Read a line from the file and check whether the third argument is
numerical'''
    line_ok = False
    while not line_ok:
        line = f.readline()
        if line == "":
            return ""
        li = line.split()
        if li[2].isdigit():
            line_ok = True
    li[2] = int(li[2])
    li[3] = int(li[3])
    return li

def main():
    list1 = [];
    if len(sys.argv) < 3:
        print "Use : get_overlap <file1> <file2>"
        os._exit(1)
    file1 = open(sys.argv[1])
    file2 = open(sys.argv[2])

    list1.append(read_line(file1))
    line2 = read_line(file2)
    while list1[-1] != "" and line2 != "":
        if line2[3] <= list1[0][2]:
            # end element 2 smaller then start smallest element 1
            # read new element 2
            line2 = read_line(file2)
        elif line2[2] > list1[-1][3]:
            # start element 2 bigger than end last element 1
            # read new element 1, remove all existing elements 1
            list1 = []
            list1.append(read_line(file1))
        else:
            # There is some overlap
            # check all elements in list in for overlap
            for line1 in list1:
                if (line2[3] >= line1[2] and line2[3] <= line1[3]) or (line1[3]
>= line2[2] and line1[3] <= line2[3]):
                    # Overlap
                    print "Overlappping elements :\n", line1, "\n", line2
            line2 = read_line(file2)
    file1.close()
    file2.close()


if __name__ == '__main__':
    main()
```
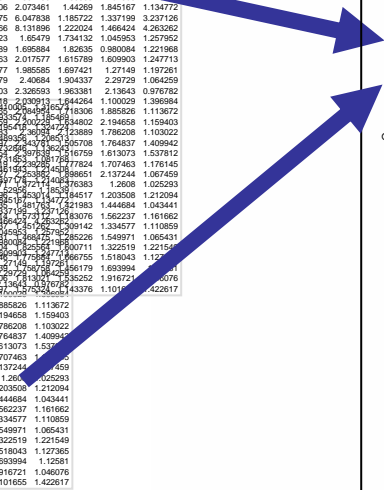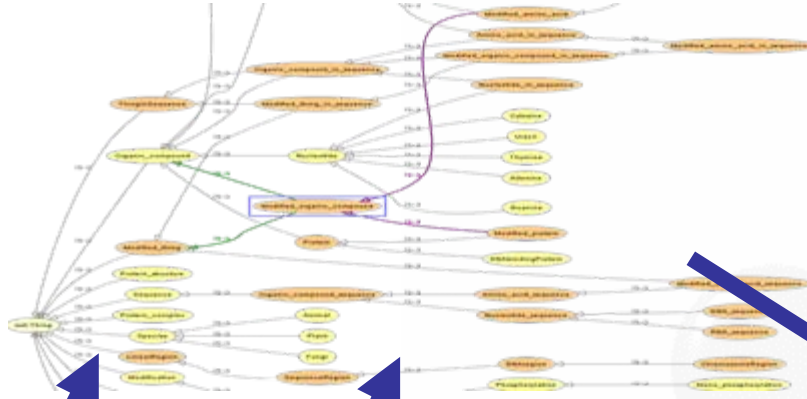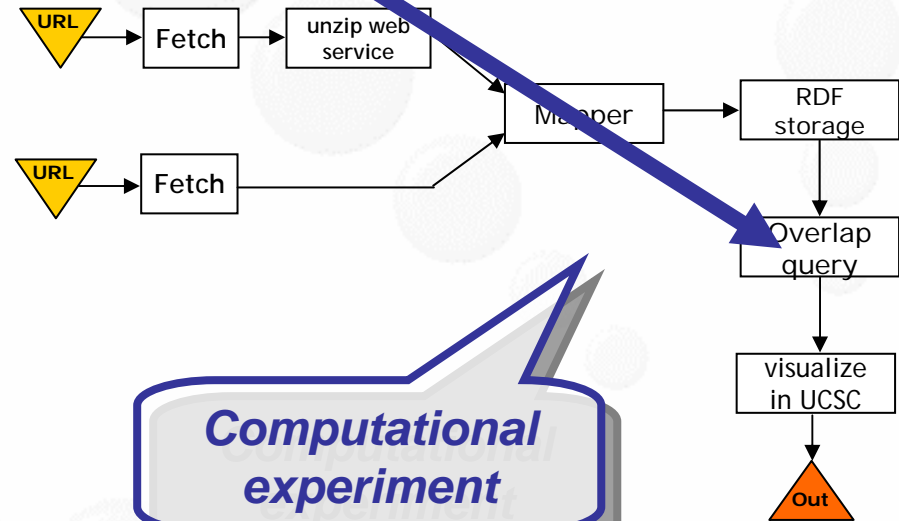
# HistOn (Histone Ontology)



*Computational experiment*

Please visit
http://integrativebioinformatics.nl
/semanticdataintegration.html

vl·e

**Our concept**
Manuscript in preparation
virtual laboratory for e·science

IBU MAD

```
SELECT * FROM {x}
        myModel:chromosomeIdentifier {y}
        USING NAMESPACE
        myModel = <http://www.somewhere.org/myModel.owl#>
```

**Table 1: Sesame query results (also available in RDF)**

| x | y |
|---|---|
| http://www.encodeSangerChipH3K4me3.xml/row1 | "chr1" |
| http://www.encodeSangerChipH3K4me3.xml/row2 | "chr1" |

2 results found in 1 ms.

# The histone code case

*'Elucidate the relationship between the Histone code, DNA sequence, and transcriptional activity'*

# The histone code case

**Multifaceted problem, no 1:1 relationships**
**Requires explorative analysis (phenomenon discovery)**
**Requires 'large-scale' data integration**

$$\Downarrow$$

**Semantic modelling approach?**

# The histone code case

**Multifaceted
no 1:1 relationships
Requires data integration and modelling**

$$\Downarrow$$

**Semantic modelling approach?**