

ACGTCAGCTAGCTAGCTAGTCAGTCGATCGATGTGTACATGCAATCGTAGCTAGCTAGCTAGCT  
TGCATGCAACTGATGCATGCATGCTGATCGATGCATGCT ROSARIO.PIRO@UNITO.IT CTAGC

# Motivation

01010101011110100101001010101000101011110000101101001010111001101011

- Linkage analysis can help to identify disease-associated loci, but:

ACGTCAGCTAGCTAGCTAGTCAGTCGATCGATGTGTACATGCAATCGTAGCTAGCTAGCTAGCT  
TGCATGCAACTGATGCATGCATGCTGATCGATGCATGCT ROSARIO.PIRO@UNITO.IT CTAGC

01010101011110100101001010101000101011110000101101001010111001101011

- Linkage analysis can help to identify disease-associated loci, but:
  - often hundreds of *positional candidate genes*

ACGTCAGCTAGCTAGCTAGTCAGTCGATCGATGTGTACATGCAATCGTAGCTAGCTAGCTAGCT  
TGCATGCAACTGATGCATGCATGCTGATCGATGCATGCT ROSARIO.PIRO@UNITO.IT CTAGC

01010101011110100101001010101000101011110000101101001010111001101011

- Linkage analysis can help to identify disease-associated loci, but:
  - often hundreds of *positional candidate genes*
  - Online Mendelian Inheritance in Man (OMIM): currently over 1000 “orphan” disease loci for phenotypes with **unknown molecular basis** but mapped locus

- Linkage analysis can help to identify disease-associated loci, but:
  - often hundreds of *positional candidate genes*
  - Online Mendelian Inheritance in Man (OMIM): currently over 1000 “orphan” disease loci for phenotypes with **unknown molecular basis** but mapped locus
  - next generation sequencing techniques?
    - large-scale re-sequencing of Chr. X in 208 XLMR families indicated 3 novel (and several known) disease genes, but in most cases the genetic cause could not be reliably determined, although many mutations (even truncating) were found.  
(Tarpey et al., Nat. Genet., 2009)

- => Computational ***disease gene prioritization***

01010101011110100101001010101000101011110000101101001010111001101011

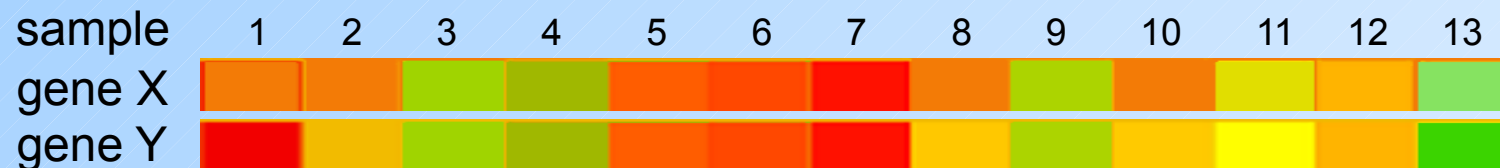
- | sample | 1   | 2   | 3   | 4   | 5   | 6   | 7   | 8   | 9   | 10  | 11  | 12  | 13  |
|--------|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|
| gene X | 0.8 | 0.8 | 0.9 | 0.7 | 0.8 | 0.8 | 0.9 | 0.8 | 0.9 | 0.8 | 0.9 | 0.8 | 0.9 |

ACGTCAGCTAGCTAGCTAGTCAGTCGATCGATGTGTACATGCAATCGTAGCTAGCTAGCTAGCT  
TGCATGCAACTGATGCATGCATGCTGATCGATGCATGCT ROSARIO.PIRO@UNITO.IT CTAGC



01010101011110100101001010101000101011110000101101001010111001101011

- “Traditional” high-throughput expression data



ACGTCAGCTAGCTAGCTAGTCAGTCGATCGATGTGTACATGCAATCGTAGCTAGCTAGCTAGCT  
TGCATGCAACTGATGCATGCATGCTGATCGATGCATGCT ROSARIO.PIRO@UNITO.IT CTAGC



01010101011110100101001010101000101011110000101101001010111001101011

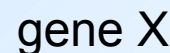
- | sample | 1      | 2      | 3           | 4           | 5      | 6      | 7   | 8      | 9           | 10     | 11     | 12     | 13          |
|--------|--------|--------|-------------|-------------|--------|--------|-----|--------|-------------|--------|--------|--------|-------------|
| gene X | orange | orange | light green | light green | orange | orange | red | orange | light green | orange | yellow | orange | light green |
| gene Y | red    | yellow | light green | light green | orange | orange | red | yellow | light green | yellow | yellow | orange | green       |

- carry **limited spatial information**, no precise 3D localization
- have a **sparse coverage** (arbitrary, not-consecutive positions)
- problem for tissues/organs with a high degree of spatial organization, e.g. the central nervous system (CNS)

01010101011110100101001010101000101011110000101101001010111001101011

- | sample | 1      | 2      | 3           | 4           | 5      | 6      | 7   | 8      | 9           | 10     | 11     | 12     | 13          |
|--------|--------|--------|-------------|-------------|--------|--------|-----|--------|-------------|--------|--------|--------|-------------|
| gene X | orange | orange | light green | light green | orange | orange | red | orange | light green | orange | yellow | orange | light green |
| gene Y | red    | yellow | light green | light green | orange | orange | red | yellow | light green | yellow | yellow | orange | green       |

- © Allen Institute



ACGTCAGCTAGCTAGCTAGTCAGTCGATCGATGTGTACATGCAATCGTAGCTAGCTAGCTAGCT  
TGCATGCAACTGATGCATGCATGCTGATCGATGCATGCT ROSARIO.PIRO@UNITO.IT CTAGC

01010101011110100101001010101000101011110000101101001010111001101011

- | sample | 1      | 2      | 3           | 4           | 5      | 6      | 7   | 8      | 9           | 10     | 11     | 12     | 13          |
|--------|--------|--------|-------------|-------------|--------|--------|-----|--------|-------------|--------|--------|--------|-------------|
| gene X | orange | orange | light green | light green | orange | orange | red | orange | light green | orange | yellow | orange | light green |
| gene Y | red    | yellow | light green | light green | orange | orange | red | yellow | light green | yellow | yellow | orange | green       |

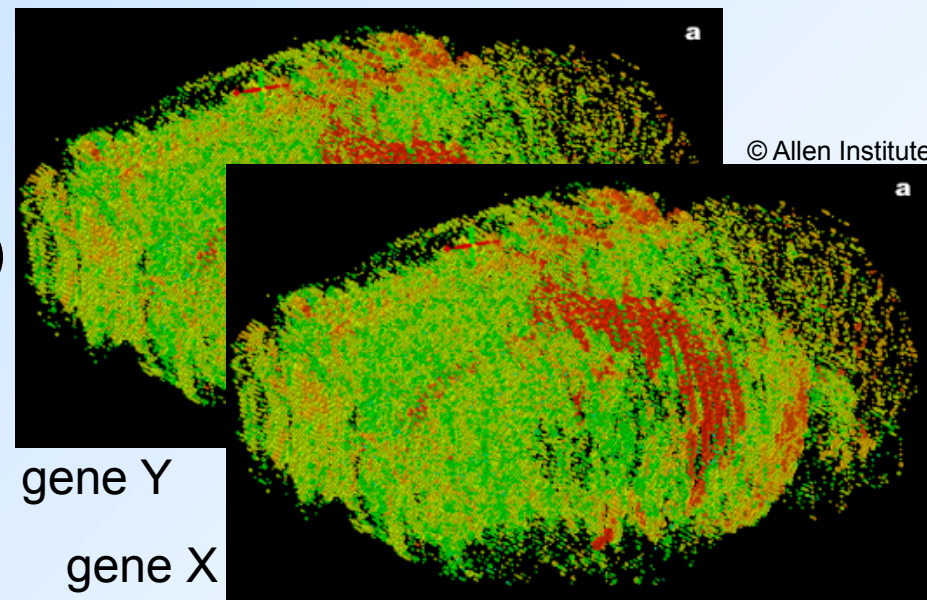
- 
- gene Y
- gene X
- a
- © Allen Institute

01010101011110100101001010101000101011110000101101001010111001101011

- | sample | 1      | 2      | 3           | 4           | 5      | 6      | 7   | 8      | 9           | 10     | 11     | 12     | 13          |
|--------|--------|--------|-------------|-------------|--------|--------|-----|--------|-------------|--------|--------|--------|-------------|
| gene X | orange | orange | light green | light green | orange | orange | red | orange | light green | orange | yellow | orange | light green |
| gene Y | red    | yellow | light green | light green | orange | orange | red | yellow | light green | yellow | yellow | orange | green       |

- Allen (Mouse) Brain Atlas

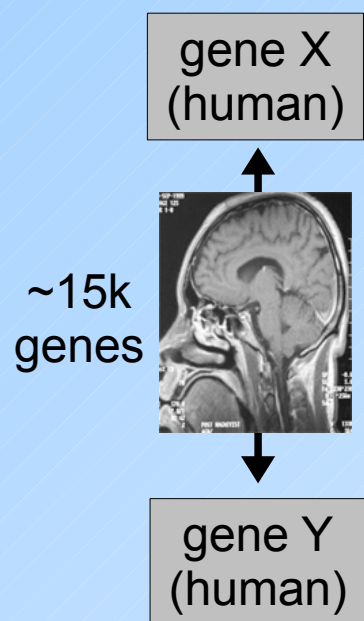
- <http://www.brain-map.org/>
- ~20k genes (~18k with Entrez ID)
- *in situ* hybridization (ISH)
- covers the entire mouse brain
- expression levels for “voxels” (cubes) of 200  $\mu\text{m}$  side length





01010101011110100101001010101000101011110000101101001010111001101011

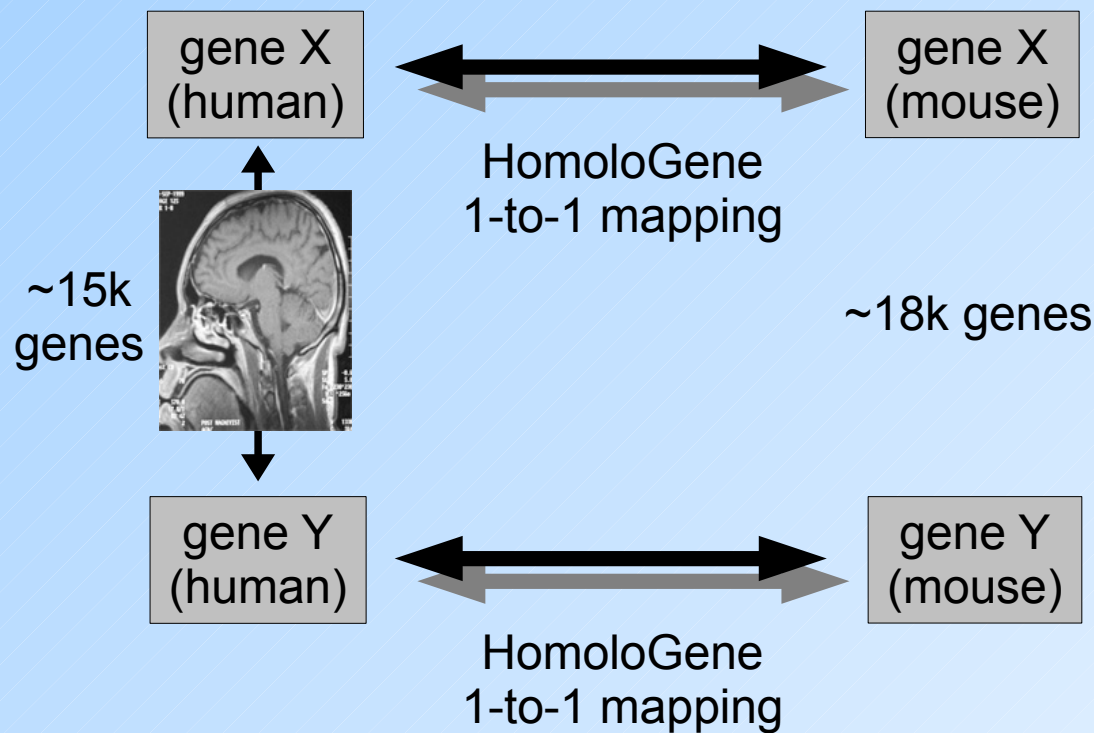
- We use mouse expression data also for **human Mendelian disorders**



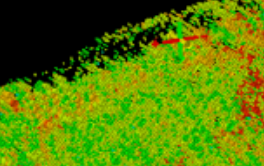
ACGTCAGCTAGCTAGCTAGTCAGTCGATCGATGTGTACATGCAATCGTAGCTAGCTAGCTAGCT  
TGCATGCAACTGATGCATGCATGCTGATCGATGCATGCT ROSARIO.PIRO@UNITO.IT CTAGC

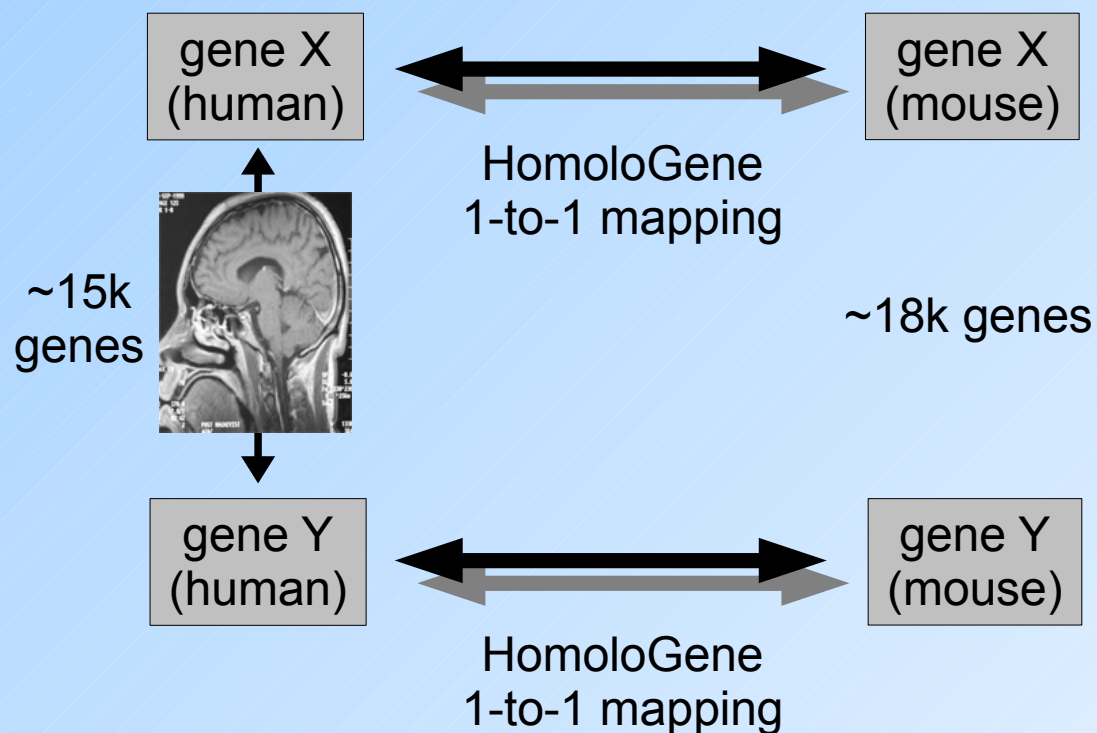
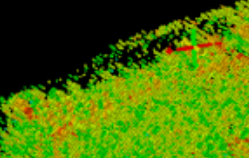
010101010111110100101001010101000101011110000101101001010111001101011

- We use mouse expression data also for **human Mendelian disorders**



01010101011110100101001010101000101011110000101101001010111001101011

- 





01010101011110100101001010101000101011110000101101001010111001101011

- *Step 1 – select “reference genes”*
  - **genes already known** to be involved in the given phenotype or in similar phenotypes (MimMiner; *van Driel et al., Eur. J. Hum. Genet., 2006*)

01010101011110100101001010101000101011110000101101001010111001101011

- *Step 1 – select “reference genes”*
  - **genes already known** to be involved in the given phenotype or in similar phenotypes (MimMiner; *van Driel et al., Eur. J. Hum. Genet., 2006*)
    - the method can be applied also to OMIM phenotypes of so far unknown molecular basis

- *Step 1 – select “reference genes”*
  - **genes already known** to be involved in the given phenotype or in similar phenotypes (MimMiner; *van Driel et al., Eur. J. Hum. Genet., 2006*)
    - the method can be applied also to OMIM phenotypes of so far unknown molecular basis
- *Step 2 – for each reference gene  $g$ :*
  - build **ranked co-expression lists** by ranking all other genes according to their (Pearson) correlation with  $g$

01010101011110100101001010101000101011110000101101001010111001101011

- ACGTCAGCTAGCTAGCTAGTCAGTCGATCGATGTGTACATGCAATCGTAGCTAGCTAGCTAGCT  
TGCATGCAACTGATGCATGCATGCTGATCGATGCATGCT ROSARIO.PIRO@UNITO.IT CTAGC

# Prioritization approach (1)

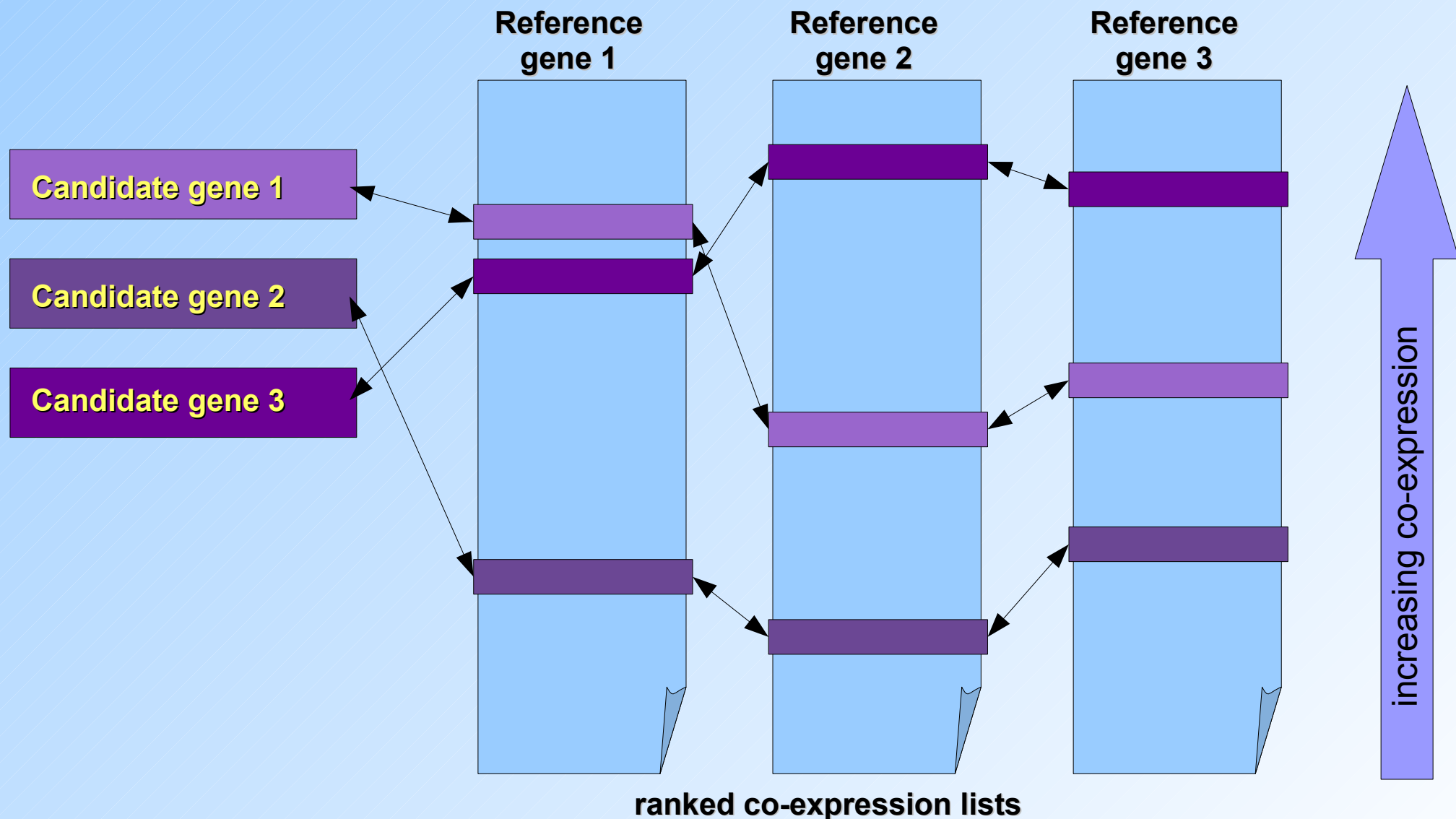
01010101011110100101001010101000101011110000101101001010111001101011

- *Step 1 – select “reference genes”*
  - **genes already known** to be involved in the given phenotype or in similar phenotypes (MimMiner; *van Driel et al., Eur. J. Hum. Genet., 2006*)
    - the method can be applied also to OMIM phenotypes of so far unknown molecular basis
- *Step 2 – for each reference gene  $g$ :*
  - build **ranked co-expression lists** by ranking all other genes according to their (Pearson) correlation with  $g$
- *Step 3 – for each positional candidate  $c$ :*
  - determine the **candidate's (relative) ranks** in the ranked co-expression lists of all reference genes
- *Step 4 – **score the candidates** (rank product)*

ACGTCAGCTAGCTAGCTAGTCAGTCGATCGATGTGTACATGCAATCGTAGCTAGCTAGCTAGCT  
TGCATGCAACTGATGCATGCATGCTGATCGATGCATGCT ROSARIO.PIRO@UNITO.IT CTAGC

# Prioritization approach (2)

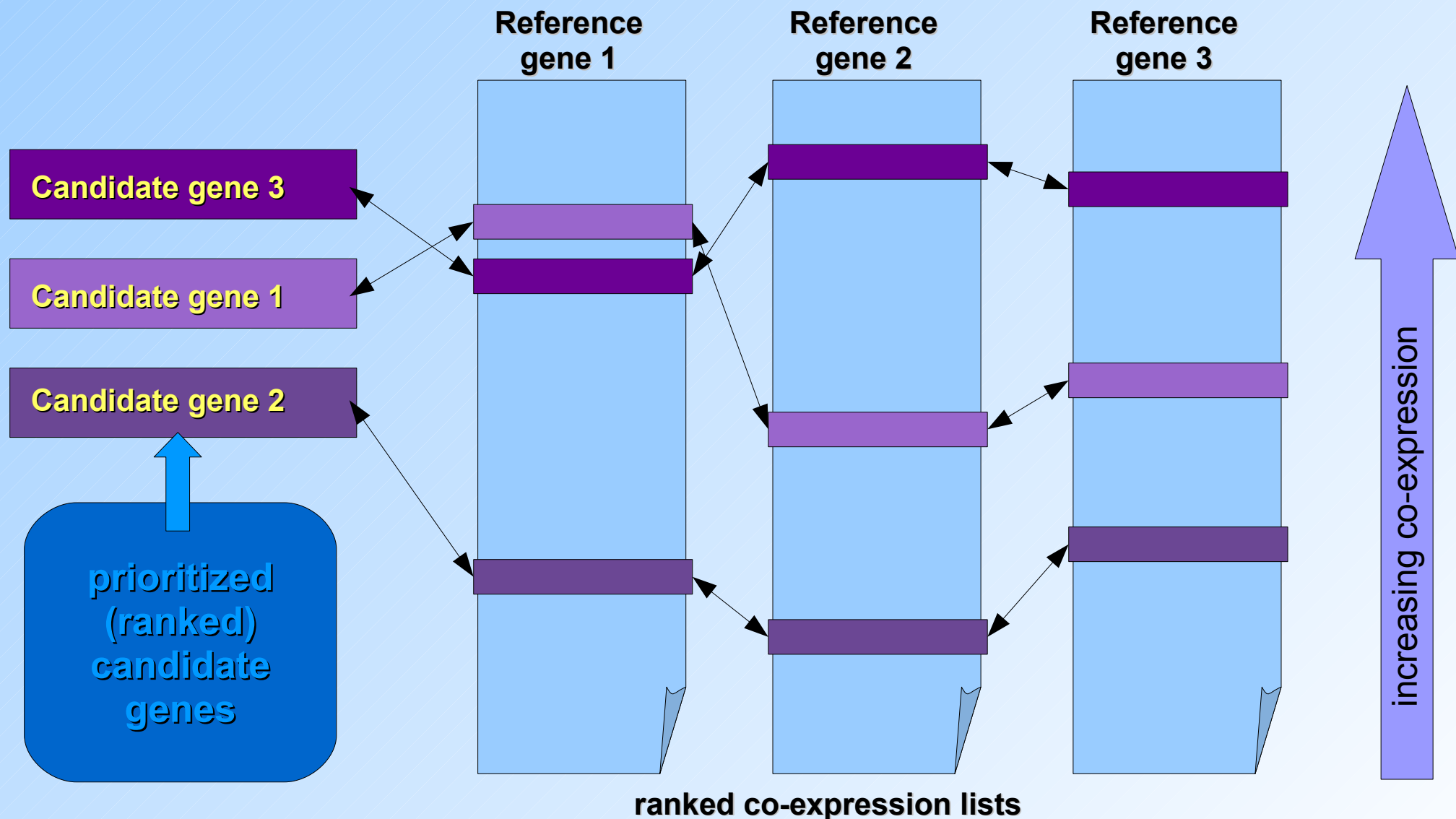
010101010111110100101001010101000101011110000101101001010111001101011



ACGTCAGCTAGCTAGCTAGCTCAGTCGATCGATGTGTACATGCAATCGTAGCTAGCTAGCTAGCT  
TGCATGCAACTGATGCATGCATGCTGATCGATGCATGCT ROSARIO.PIRO@UNITO.IT CTAGC



01010101011110100101001010101000101011110000101101001010111001101011



ACGTCAGCTAGCTAGCTAGTCAGTCGATCGATGTGTACATGCAATCGTAGCTAGCTAGCTAGCT  
TGCATGCAACTGATGCATGCATGCTGATCGATGCATGCT ROSARIO.PIRO@UNITO.IT CTAGC



01010101011110100101001010101000101011110000101101001010111001101011

- Large-scale leave-one-out validations for ...
  - ... **known CNS-related gene-phenotype associations** from different databases
    - Mouse Genome Database (MGD)
    - Online Mendelian Inheritance in Man (OMIM)
  - ... **different sizes of loci**

- ACGTCAGCTAGCTAGCTAGTCAGTCGATCGATGTGTACATGCAATCGTAGCTAGCTAGCTAGCT  
TGCATGCAACTGATGCATGCATGCTGATCGATGCATGCT ROSARIO.PIRO@UNITO.IT CTAGC

01010101011110100101001010101000101011110000101101001010111001101011

- (Céciz et al., Trends Genet., 2009)

01010101011110100101001010101000101011110000101101001010111001101011

- ACGTCAGCTAGCTAGCTAGTCAGTCGATCGATGTGTACATGCAATCGTAGCTAGCTAGCTAGCT  
TGCATGCAACTGATGCATGCATGCTGATCGATGCATGCT ROSARIO.PIRO@UNITO.IT CTAGC

01010101011110100101001010101000101011110000101101001010111001101011

- Challenging, complex disorder:  
~90 genes on Chr. X are known to be involved in  
some form of intellectual disability  
(Céciz et al., Trends Genet., 2009)
- A similar number probably remains to be identified  
(Céciz et al., Trends Genet., 2009)
- Recent re-sequencing study of the exome of Chr. X  
(Tarpey et al., Nat. Genet., 2009)
  - 208 affected families
  - mutations in 3 novel (and several known) XLMR genes
  - many other mutations (also truncating), but in for most of the  
cases the genetic cause could not be reliably determined

01010101011110100101001010101000101011110000101101001010111001101011

- Evaluation:
  - candidate genes: all re-sequenced genes
  - reference genes: only genes known to be involved in similar phenotypes (none of the candidates from Chr. X)
  - can we “predict” known disease genes **pretending XLMR to be a phenotype of unknown molecular basis?**



01010101011110100101001010101000101011110000101101001010111001101011

- Evaluation:
  - candidate genes: all re-sequenced genes
  - reference genes: only genes known to be involved in similar phenotypes (none of the candidates from Chr. X)
  - can we “predict” known disease genes **pretending XLMR to be a phenotype of unknown molecular basis?**
- Prediction:
  - reference genes: XLMR genes
  - candidate genes: all other re-sequenced genes
  - can we find **promising novel candidates** for XLMR?

ACGTCAGCTAGCTAGCTAGTCAGTCGATCGATGTGTACATGCAATCGTAGCTAGCTAGCTAGCT  
TGCATGCAACTGATGCATGCATGCTGATCGATGCATGCT ROSARIO.PIRO@UNITO.IT CTAGC



01010101011110100101001010101000101011110000101101001010111001101011

- Evaluation:
  - candidate genes: all re-sequenced genes
  - reference genes: only genes known to be involved in similar phenotypes (none of the candidates from Chr. X)
  - can we “predict” known disease genes **pretending XLMR to be a phenotype of unknown molecular basis?**
- Prediction:
  - reference genes: XLMR genes
  - candidate genes: all other re-sequenced genes
  - can we find **promising novel candidates** for XLMR?
- Important: evaluation and prediction have completely distinct sets of reference genes!

ACGTCAGCTAGCTAGCTAGTCAGTCGATCGATGTGTACATGCAATCGTAGCTAGCTAGCTAGCT  
TGCATGCAACTGATGCATGCATGCTGATCGATGCATGCT ROSARIO.PIRO@UNITO.IT CTAGC

01010101011110100101001010101000101011110000101101001010111001101011

1	BRWD3
2	IRAK1
3	SYP
4	BIRC4
5	MAGED1
6	MORF4L2
7	ZNF280C
8	SYN1
9	CXorf6
10	ATP6AP2
11	HCFC1
12	PJA1
13	NGFRAP1
14	FAM50A
15	HUWE1
16	GRIA3
17	PIGA
18	OGT
19	GNL3L
20	WDR40C

ACGTCAGCTAGCTAGCTAGTCAGTCGATCGATGTGTACATGCAATCGTAGCTAGCTAGCTAGCT  
TGCATGCAACTGATGCATGCATGCTGATCGATGCATGCT ROSARIO.PIRO@UNITO.IT CTAGC

01010101011110100101001010101000101011110000101101001010111001101011

1	BRWD3
2	IRAK1
3	SYP
4	BIRC4
5	MAGED1
6	MORF4L2
7	ZNF280C
8	SYN1
9	CXorf6
10	ATP6AP2
11	HCFC1
12	PJA1
13	NGFRAP1
14	FAM50A
15	HUWE1
16	GRIA3
17	PIGA
18	OGT
19	GNL3L
20	WDR40C

- successful “re-discovery” of known **XLMR genes**

ACGTCAGCTAGCTAGCTAGTCAGTCGATCGATGTGTACATGCAATCGTAGCTAGCTAGCTAGCT  
TGCATGCAACTGATGCATGCATGCTGATCGATGCATGCT ROSARIO.PIRO@UNITO.IT CTAGC

010101010111110100101001010101000101011110000101101001010111001101011

1	BRWD3
2	IRAK1
3	SYP
4	BIRC4
5	MAGED1
6	MORF4L2
7	ZNF280C
8	SYN1
9	CXorf6
10	ATP6AP2
11	HCFC1
12	PJA1
13	NGFRAP1
14	FAM50A
15	HUWE1
16	GRIA3
17	PIGA
18	OGT
19	GNL3L
20	WDR40C

- successful “re-discovery” of known **XLMR genes**

1	MORF4L2
2	PJA1
3	ZNF280C
4	MAGED1
5	MAGEE1
6	BIRC4
7	GRIPAP1
8	CXorf6
9	GNL3L
10	FAM50A
11	PGRMC1
12	GPM6B
13	IRAK1
14	HCFC1
15	PIGA
16	RPS4X
17	REPS2
18	ARMCX2
19	DRP2
20	MED14

ACGTCAGCTAGCTAGCTAGTCAGTCGATCGATGTGTACATGCAATCGTAGCTAGCTAGCTAGCTA  
TGCATGCAACTGATGCATGCATGCTGATCGATGCATGCT ROSARIO.PIRO@UNITO.IT CTAGC

# Case study: XLMR - Results

010101010111110100101001010101000101011110000101101001010111001101011

# Evaluation

2	IRAK1
4	BIRC4
5	MAGED1
6	MORF4L2
7	ZNF280C
9	CXorf6
11	HCFC1
12	PJA1
13	NGFRAP1
14	FAM50A
17	PIGA
18	OGT
19	GNL3L
20	WDR40C

- successful “re-discovery” of known **XLMR genes**

# Prediction

1	MORF4L2
2	PJA1
3	ZNF280C
4	MAGED1
5	MAGEE1
6	BIRC4
7	GRIPAP1
8	CXorf6
9	GNL3L
10	FAM50A
11	PGRMC1
12	GPM6B
13	IRAK1
14	HCFC1
15	PIGA
16	RPS4X
17	REPS2
18	ARMCX2
19	DRP2
20	MED14

ACGTCAGCTAGCTAGCTAGTCAGTCGATCGATGTGTACATGCAATCGTAGCTAGCTAGCTAGCTA  
TGCATGCAACTGATGCATGCATGCTGATCGATGCATGCT ROSARIO.PIRO@UNITO.IT CTAGC

# Case study: XLMR - Results

01010101011110100101001010101000101011110000101101001010111001101011

# Evaluation

2	IRAK1
4	BIRC4
5	MAGED1
6	MORF4L2
7	ZNF280C
9	CXorf6
11	HCFC1
12	PJA1
13	NGFRAP1
14	FAM50A
17	PIGA
18	OGT
19	GNL3L
20	WDR40C

- successful “re-discovery” of known **XLMR genes**
- very strong **overlap** between evaluation and prediction ( $P=6.62 \times 10^{-14}$ )

# Prediction

1	MORF4L2
2	PJA1
3	ZNF280C
4	MAGED1
5	MAGEE1
6	BIRC4
7	GRIPAP1
8	CXorf6
9	GNL3L
10	FAM50A
11	PGRMC1
12	GPM6B
13	IRAK1
14	HCFC1
15	PIGA
16	RPS4X
17	REPS2
18	ARMCX2
19	DRP2
20	MED14

ACGTCAGCTAGCTAGCTAGTCAGTCGATCGATGTGTACATGCAATCGTAGCTAGCTAGCTAGCTA  
TGCATGCAACTGATGCATGCATGCTGATCGATGCATGCT ROSARIO.PIRO@UNITO.IT CTAGC

01010101011110100101001010101000101011110000101101001010111001101011

# Prediction

ACGTCAGCTAGCTAGCTAGTCAGTCGATCGATGTGTACATGCAATCGTAGCTAGCTAGCTAGCTA  
TGCATGCAACTGATGCATGCATGCTGATCGATGCATGCT ROSARIO.PIRO@UNITO.IT CTAGC



- Proof-of-concept study
  - no extensive comparison with other data sources
  - no optimization
    - alternative scoring functions instead of the rank product?

- Proof-of-concept study
  - no extensive comparison with other data sources
  - no optimization
    - alternative scoring functions instead of the rank product?
- Spatial information not fully exploited
  - alternative measures for the similarity of expression profiles?  
use HRC instead of the Pearson correlation coefficient?  
(HRC = histogram-row-column; *Liu et al., BMC Sys. Biol., 2007*)

01010101011110100101001010101000101011110000101101001010111001101011

- **Proof-of-concept study**
  - no extensive comparison with other data sources
  - no optimization
    - alternative scoring functions instead of the rank product?
- **Spatial information not fully exploited**
  - alternative measures for the similarity of expression profiles?  
use HRC instead of the Pearson correlation coefficient?  
(HRC = histogram-row-column; *Liu et al., BMC Sys. Biol., 2007*)
- **Spatial 3D expression data remain an exception**
  - Human Brain Atlas (Allen Institute) to be completed in 2013
    - our work can be considered a pioneer study towards a direct application of human 3D expression data
  - so far only “normal” brain tissue; some potentially interesting applications would require disease conditions
    - “differential spatial expression”?

ACGTCAGCTAGCTAGCTAGTCAGTCGATCGATGTGTACATGCAATCGTAGCTAGCTAGCTAGCT  
TGCATGCAACTGATGCATGCATGCTGATCGATGCATGCT ROSARIO.PIRO@UNITO.IT CTAGC

01010101011110100101001010101000101011110000101101001010111001101011

- We have shown that spatially mapped (3D) gene expression data can be successfully exploited for candidate gene prioritization
  - mouse CNS-related phenotypes
  - human CNS-related Mendelian disorders

01010101011110100101001010101000101011110000101101001010111001101011

- ACGTCAGCTAGCTAGCTAGTCAGTCGATCGATGTGTACATGCAATCGTAGCTAGCTAGCTAGCT  
TGCATGCAACTGATGCATGCATGCTGATCGATGCATGCT ROSARIO.PIRO@UNITO.IT CTAGC

01010101011110100101001010101000101011110000101101001010111001101011

- Nijmegen Centre for  
Molecular Life Sciences  
(The Netherlands)



ACGTCAGCTAGCTAGCTAGTCAGTCGATCGATGTGTACATGCAATCGTAGCTAGCTAGCTAGCT  
TGCATGCAACTGATGCATGCATGCTGATCGATGCATGCT ROSARIO.PIRO@UNITO.IT CTAGC

01010101011110100101001010101000101011110000101101001010111001101011



ACGTCAGCTAGCTAGCTAGTCAGTCGATCGATGTGTACATGCAATCGTAGCTAGCTAGCTAGCT  
TGCATGCAACTGATGCATGCATGCTGATCGATGCATGCT ROSARIO.PIRO@UNITO.IT CTAGC

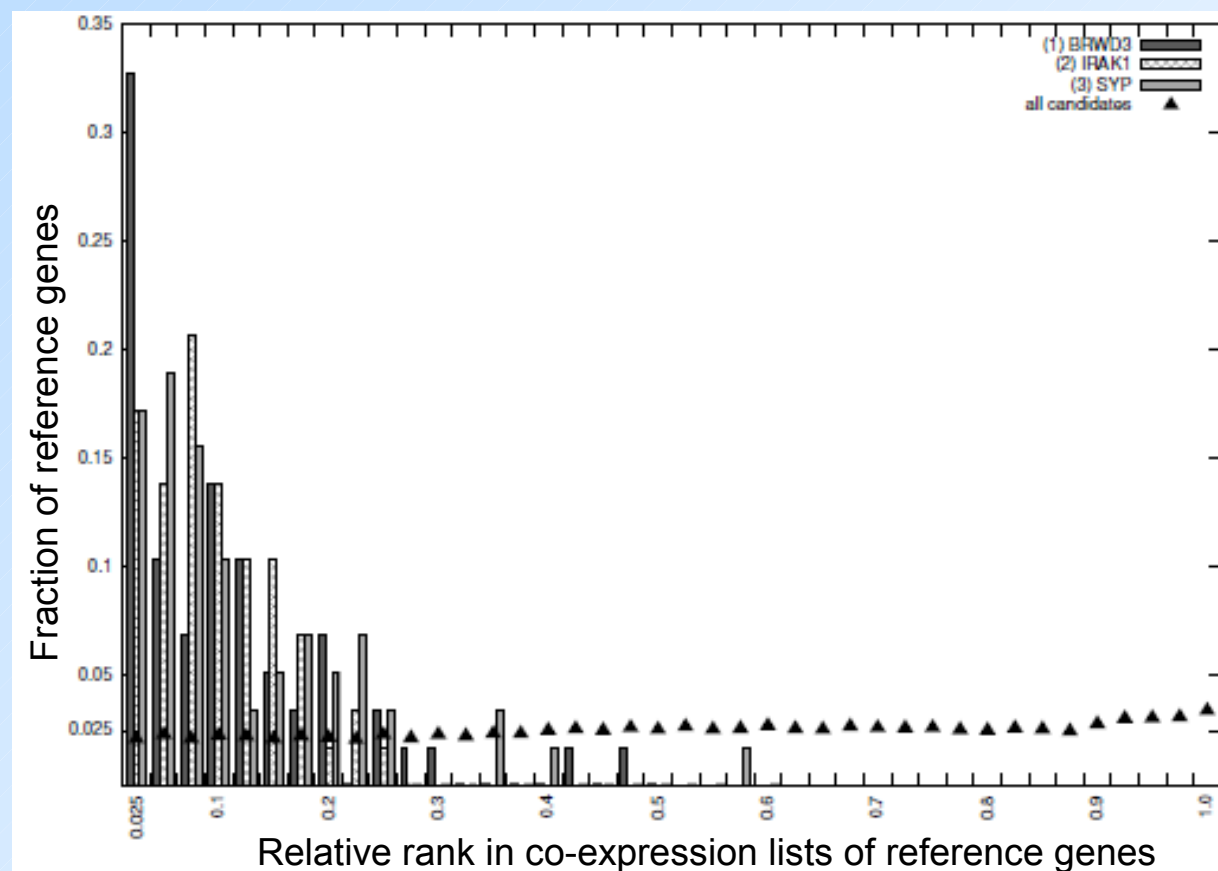
# Case study: XLMR - Results

01010101011110100101001010101000101011110000101101001010111001101011

# Evaluation

1	BRWD3
2	IRAK1
3	SYP
4	BIRC4
5	MAGED1
6	MORF4L2
7	ZNF280C
8	SYN1
9	CXorf6
10	ATP6AP2
11	HCFC1
12	PJA1
13	NGFRAP1
14	FAM50A
15	HUWE1
16	GRIA3
17	PIGA
18	OGT
19	GNL3L
20	WDR40C

- successful “re-discovery” of known **XLMR genes**



ACGTCAGCTAGCTAGCTAGTCAGTCGATCGATGTGTACATGCAATCGTAGCTAGCTAGCTAGCT  
TGCATGCAACTGATGCATGCATGCTGATCGATGCATGCT ROSARIO.PIRO@UNITO.IT CTAGC