

# Experimental Design

Mikhail Dozmorov  
Fall 2017

## What is experimental design?

The organization of an experiment, to ensure that *the right type of data*, and *enough of it*, is available to answer the questions of interest as clearly and efficiently as possible.

## What characterizes well-designed experiments?

- Effects can be estimated unambiguously and without bias.
- Estimates are precise.
- Protected from possible one-off events that might compromise the results.
- Easy to conduct.
- Easy to analyse and interpret.
- Maximum information obtained for fixed time, resources, and samples.
- Applicability of the findings to a wide variety of subjects, conditions, and situations.

3/23

## Why Design an Experiment?

- The goal of an experiment dictates everything from how the samples are collected to how the data are generated
- The design of the analytical protocol should be reflected in the design
  - Do we have enough replicates?
  - Do we have sufficient controls?
  - Do we collect samples and data to avoid confounding and batch effects?

4/23

# Types of Experiments

## Class Comparison

- Can I find genes that distinguish between two classes, such as tumor and normal?

## Class Discovery

- Given what I think is a uniform group of samples, can I find subsets that are biologically meaningful?

## Classification

- Given a set of samples in different classes, can I assign a new, unknown sample to one of the classes?

## Large-scale Functional Studies

- Can I discover a causative mechanism associated with the distinction between classes? These are often not perfectly distinct.

5/23

# What affects the outcome of an experiment?

$$\text{Outcome} = \underbrace{\text{Treatment effects}} + \underbrace{\text{Biological effects}} + \underbrace{\text{Technical effects}} + \underbrace{\text{Error}}$$

Environment	Sex	Technician	Experimental
Compound	Age	Batch	Treatment
Inhibitor	Weight	Plate	Sampling
siRNA	Litter	Cage	Measurement
Dose	Genotype	Array	
Time	Species	Day	
	Cell line	Order	
		Source	

6/23

## What is bad Bad experimental design - examples

Treatment I

M M M M M M M M

Treatment II

F F F F F F F F

7/23

## What is bad Bad experimental design - examples

Treatment I

M M M M M M

Treatment II

F F F F F F

**Confounding!**

8/23

## What is bad Bad experimental design - examples

Analysis batch I / Study center I / Processing protocol I ...

Tr Tr Tr Tr Tr Tr Tr Tr

Analysis batch II / Study center II / Processing protocol II ...

Ctl Ctl Ctl Ctl Ctl Ctl Ctl Ctl

9/23

## What would be a better experimental design?

- Process all samples at the same time/in one batch (not always feasible)
- Minimize confounding as much as possible through
  - blocking
  - randomization
- The batch effect will still be there, but with an appropriate design we can account for it

10/23

# Principles of experimental design

- **Replication.** It allows the experimenter to obtain an estimate of the experimental error
- **Randomization.** It requires the experimenter to use a random choice of every factor that is not of interest but might influence the outcome of the experiment. Such factors are called nuisance factors
- **Blocking.** Creating homogeneous blocks of data in which a nuisance factor is kept constant while the factor of interest is allowed to vary. Used to increase the accuracy with which the influence of the various factors is assessed in a given experiment
- **Block what you can, randomize what you cannot**

11/23

## Replicates

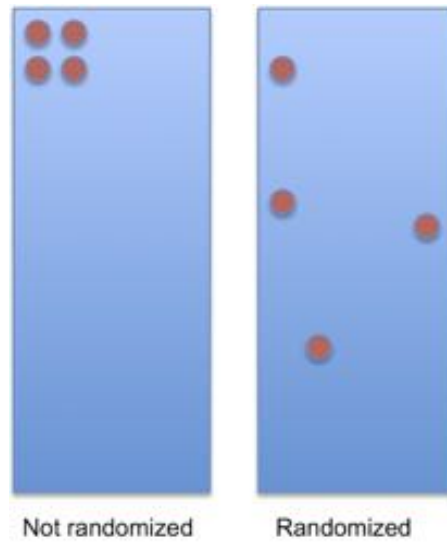


- **Technical** replicates and **Biological** replicates
- Rule of thumb: for two-fold change – use 3 replicates
- Smaller change – 5 replicates

12/23

## Randomization

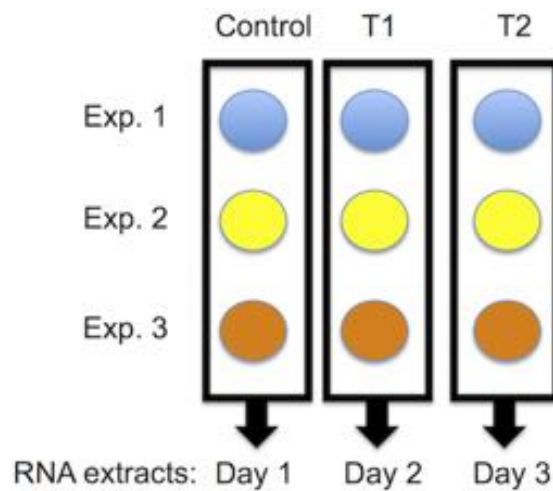
- Each gene has multiple probes – randomize their position on the slide



13/23

## Blocking

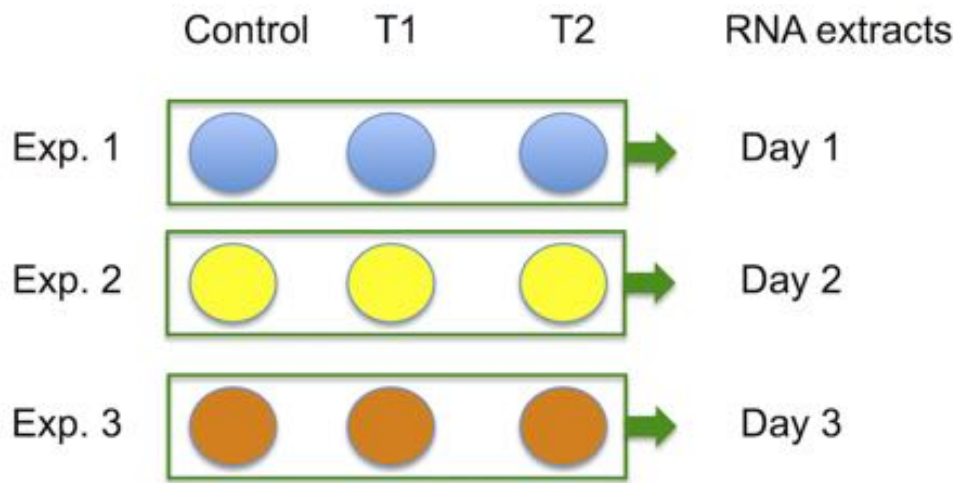
- Treatment and RNA extraction days are confounded



14/23

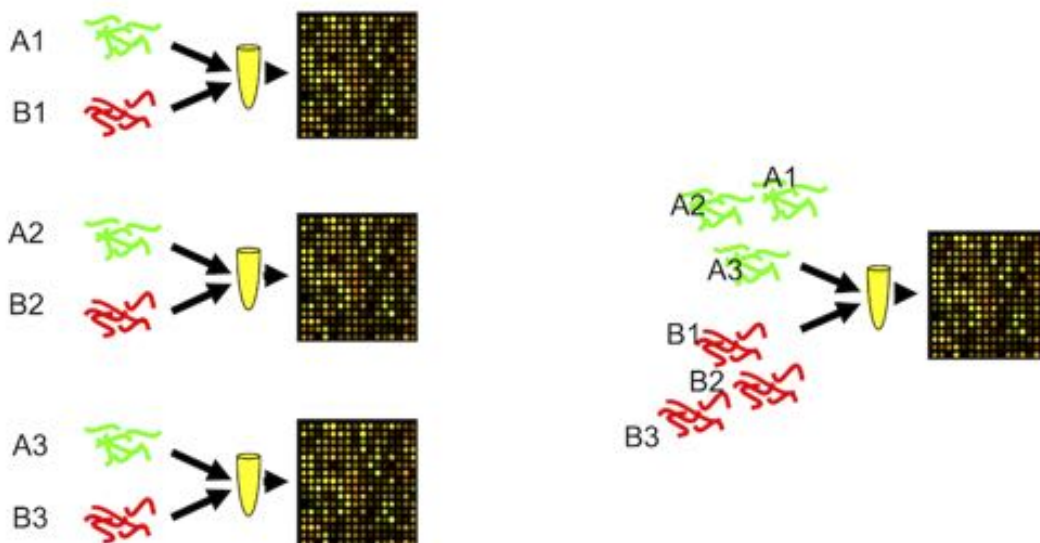
# Blocking

- Block replicated experiments



15/23

# Pooling



16/23



# Pooling

- OK when the interest is not on the individual, but on common patterns across individuals (population characteristics)
- Results in averaging -> reduces variability -> substantive features are easier to find
- Recommended when fewer than 3 arrays are used in each condition
- Beneficial when many subjects are pooled

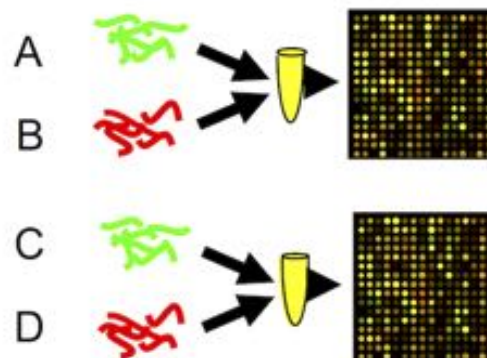
## "inference for most genes was not affected by pooling"

C. Kendziora, R. A. Irizarry, K.-S. Chen, J. D. Haag, and M. N. Gould. "On the utility of pooling biological samples in microarray experiments". PNAS March 2005, 102(12) 4252-4257

17/23

## How to allocate the samples to microarrays?

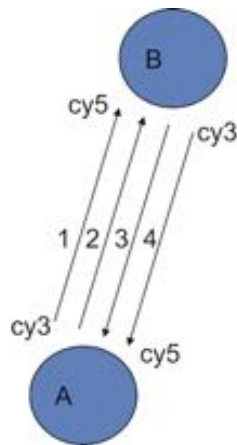
- which samples should be hybridized on the same slide?
- how different experimental designs affect outcome?
- what is the optimal design?



18/23

# Example of four-array experiment

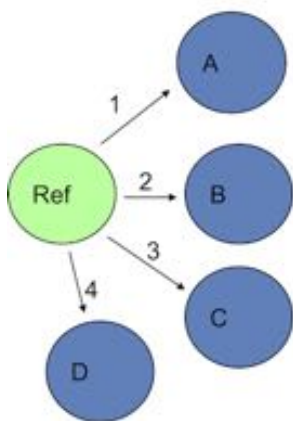
Dye swap



array	cy3	cy5	$\log(\text{cy5}/\text{cy3})$
1	A	B	$\log(B) - \log(A)$
2	A	B	$\log(B) - \log(A)$
3	B	A	$\log(A) - \log(B)$
4	B	A	$\log(A) - \log(B)$

19/23

# Common reference design

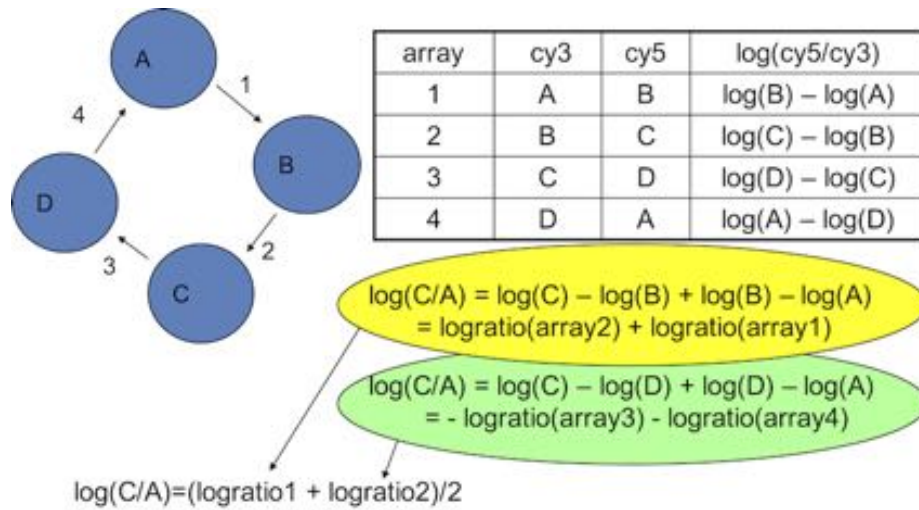


array	cy3	cy5	$\log(\text{cy5}/\text{cy3})$
1	Ref	A	$\log(A) - \log(\text{Ref})$
2	Ref	B	$\log(B) - \log(\text{Ref})$
3	Ref	C	$\log(C) - \log(\text{Ref})$
4	Ref	D	$\log(D) - \log(\text{Ref})$

$$\begin{aligned}
 \log(C/A) &= \log(C) - \log(A) \\
 &= \log(C) - \log(\text{Ref}) + \log(\text{Ref}) - \log(A) \\
 &= \log(C) - \log(\text{Ref}) - (\log(A) - \log(\text{Ref})) \\
 &= \text{logratio}(\text{array3}) - \text{logratio}(\text{array1})
 \end{aligned}$$

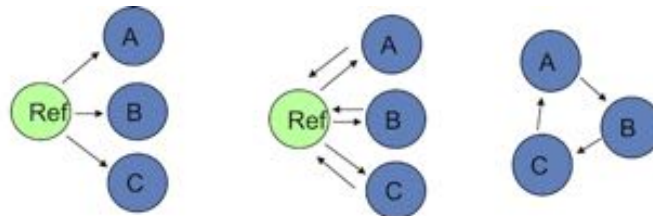
20/23

## Loop design



21/23

## Comparing the designs



	reference design	reference design with replicates	loop design
number of arrays	3	6	3
amount of RNA required per sample	1+Ref	2+Ref	2
error	2.0	1.0	0.67

22/23

## Design with all direct pairwise comparisons

