

TO COMBINE RESULTS FROM MULTIPLE STUDIES...

- We need to have some methods to combine results from the testing of the same markers = **meta-analysis methods**

TO COMBINE RESULTS FROM MULTIPLE STUDIES...

- We need to have some methods to combine results from the testing of the same markers = **meta-analysis methods**
- We need the same set of markers studied. This is not always the case (different chips used by different studies). Therefore we need methods to infer/impute the markers not typed in present study = **imputation methods**

IMPUTATIONS (AND BITS ON ANALYSIS OF IMPUTED DATA)

YURII AULCHENKO

YURII [DOT] AULCHENKO [AT] GMAIL [DOT] COM

TO COMBINE RESULTS FROM MULTIPLE STUDIES...

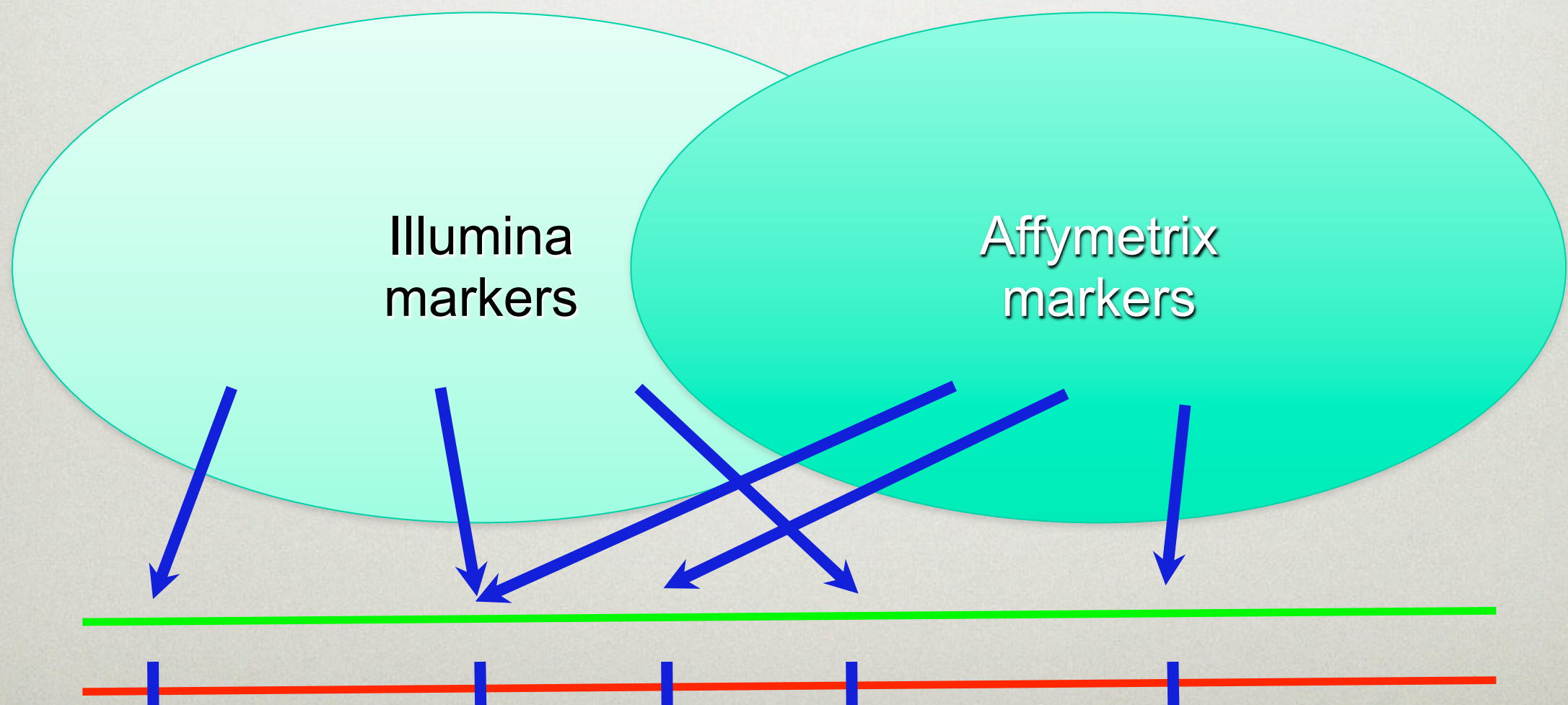
- We need to have some methods to combine results from the testing of the same markers = **meta-analysis methods**

TO COMBINE RESULTS FROM MULTIPLE STUDIES...

- We need to have some methods to combine results from the testing of the same markers = **meta-analysis methods**
- We need the same set of markers studied. This is not always the case (different chips used by different studies). Therefore we need methods to infer/impute the markers not typed in present study = **imputation methods**

HARMONIZATION THROUGH IMPUTATIONS

- Imputation to combine data across platforms
- Statistical models combining evidence across studies (meta-analysis)



Courtesy of Samuli Ripatti

IMPUTATIONS IN GENERAL

- For markers within genotyped set, based on flanking marker information
 - Predicts missing genotypes
 - Points out possible genotyping errors
- Extension: use a set of more densely typed samples (e.g. HapMap) to predict markers NOT typed in your set
- “Prediction” from imputations:
 - Not exact, but probabilistic prediction
 - E. g. some person at some SNP is predicted to have TT (90%), TG (7%), or GG (3%).

WHAT IS IMPUTATION?

Reference
haplotypes:

A	C	C	T	T	A	A	G	C	T	C	A	G	A	T	C
A	A	A	A	A	C	G	G	A	A	A	G	G	C	G	A

Own data
haplotype:

A	C	C	?	T	C	?	G	C	?	C	A	G	?	G	A
---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---

WHAT IS IMPUTATION?

Reference
haplotypes:

A	C	C	T	T	A	A	G	C	T	C	A	G	A	T	C
A	T	G	G	A	C	G	G	A	A	A	G	G	C	G	A

Own data
haplotype:

A	C	C		T	C	?	G	C	?	C	A	G	?	G	A
---	---	---	--	---	---	---	---	---	---	---	---	---	---	---	---

WHAT IS IMPUTATION?

Reference
haplotypes:

A	C	C	T	T	A	A	G	C	T	C	A	G	A	T	C
A	T	G	G	A	C	G	G	A	A	A	G	G	C	G	A

Own data
haplotype:

A	C	C	T	T	C	?	G	C	?	C	A	G	?	G	A
---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---

WHAT IS IMPUTATION?

Reference
haplotypes:

A	C	C	T	T	A	A	G	C	T	C	A	G	A	T	C
A	T	G	G	A	C	G	G	A	A	A	G	G	C	G	A

Own data
haplotype:

A	C	C	T	T	C		G	C	?	C	A	G	?	G	A
---	---	---	---	---	---	--	---	---	---	---	---	---	---	---	---

WHAT IS IMPUTATION?

Reference
haplotypes:

A	C	C	T	T	A	A	G	C	T	C	A	G	A	T	C
A	T	G	G	A	C	G	G	A	A	A	G	G	C	G	A

Own data
haplotype:

A	C	C	T	T	C	G	G	C	?	C	A	G	?	G	A
---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---

WHAT IS IMPUTATION?

Reference
haplotypes:

A	C	C	T	T	A	A	G	C	T	C	A	G	A	T	C
A	T	G	G	A	C	G	G	A	A	A	G	G	C	G	A

Own data
haplotype:

A	C	C	T	T	C	G	G	C		C	A	G	?	G	A
---	---	---	---	---	---	---	---	---	--	---	---	---	---	---	---

WHAT IS IMPUTATION?

Reference
haplotypes:

A	C	C	T	T	A	A	G	C	T	C	A	G	A	T	C
A	T	G	G	A	C	G	G	A	A	A	G	G	C	G	A

Own data
haplotype:

A	C	C	T	T	C	G	G	C	T	C	A	G	?	G	A
---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---

WHAT IS IMPUTATION?

Reference
haplotypes:

A	C	C	T	T	A	A	G	C	T	C	A	G	A	T	C
A	T	G	G	A	C	G	G	A	A	A	G	G	C	G	A

Own data
haplotype:

A	C	C	T	T	C	G	G	C	T	C	A	G		G	A
---	---	---	---	---	---	---	---	---	---	---	---	---	--	---	---

WHAT IS IMPUTATION?

Reference
haplotypes:

A	C	C	T	T	A	A	G	C	T	C	A	G	A	T	C
A	T	G	G	A	C	G	G	A	A	A	G	G	C	G	A

Own data
haplotype:

A	C	C	T	T	C	G	G	C	T	C	A	G	C	G	A
---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---

IMPUTING: GUESS THE “?”

Reference
Haplotypes

G	A	C	T	G
T	C	T	T	G
G	A	C	T	G
T	C	C	T	G
T	C	C	T	G
G	A	C	T	G
T	C	C	T	G

Sample
Genotypes

G	A	?	T	G
T	C	?	T	G

IMPUTING: GUESS THE “?”

Reference
Haplotypes

G	A	C	T	G
T	C	T	T	G
G	A	C	T	G
T	C	C	T	G
T	C	C	T	G
G	A	C	T	G
T	C	C	T	G

Sample
Genotypes

G	A	?	T	G
T	C	?	T	G

Best-guess:
C/C

Dosage:
1.75

Mixture:
C/C for 75%
C/T for 25%

IMPUTING: GUESS THE “?”

In “real life”, imputations are done using rather complex mathematical models, e.g. Hidden Markov Models, not simple guessing

Best-guess:
C/C

Dosage:
1.75

Mixture:
C/C for 75%
C/T for 25%

IMPUTED DATA

We can not tell the exact genotype, but can estimate posterior probability distribution: $P(g)=\{p_{AA},p_{AB},p_{BB}\}$

Directly typed SNPs: either AA, AB or BB. The probability distribution is *degenerate* (e.g. $\{0,1,0\}$ that is to say AB)

For imputed SNPs, the distribution is not degenerate

HOW CAN WE ANALYZE IMPUTED DATA?

- Instead of genotypes, we have probabilities that certain person at certain locus has certain genotype
- How do we relate these probabilities to an outcome (do GWAS)?

BEST GUESS

- Best guess: take the genotype with maximal posterior probability and treat it as if it was true, directly typed
- Drawback: biased estimates, reduced power

REGRESSION ONTO PROBABILITIES

Use the model

$$E[Y_i] = m + b_1 P(g_i=1) + b_2 P(g_i=2)$$

Note this is very similar to model for directly typed SNPs, with probabilities used instead of indicator variables.

Different genetic models can be formulated in the same way.

MAXIMUM LIKELIHOOD BASED ON PROBABILITIES

Define individual likelihood as

$$L_i = \text{SUM}_{g_i=\{0,1,2\}} P(g_i) P(Y_i | g_i)$$

Where

$$P(Y_i | g_i) = \text{Normal}(E[Y_i | g_i], s^2)$$

and

$$E[Y_i | g_i] = m + b_1 I(g_i=1) + b_2 I(g_i=2)$$

MAXIMUM LIKELIHOOD BASED ON PROBABILITIES

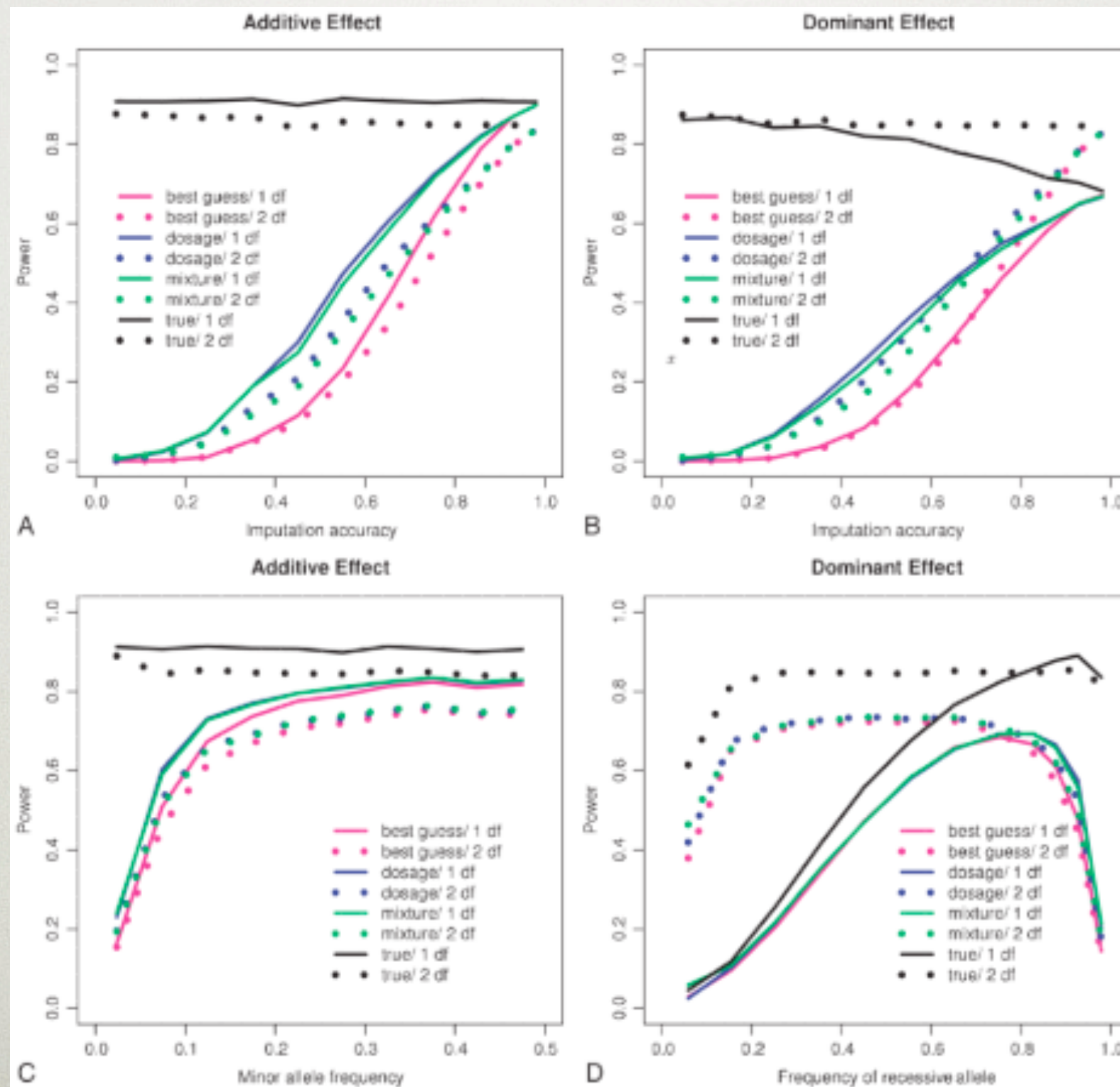
Define joint likelihood as the product of individual likelihoods

Maximize the likelihood over the parameters involved

Maximum Likelihood Ratio test can be used to test nested models and draw statistical inferences

POWER IN LARGE SAMPLES (SMALL EFFECTS)

vs.
accuracy

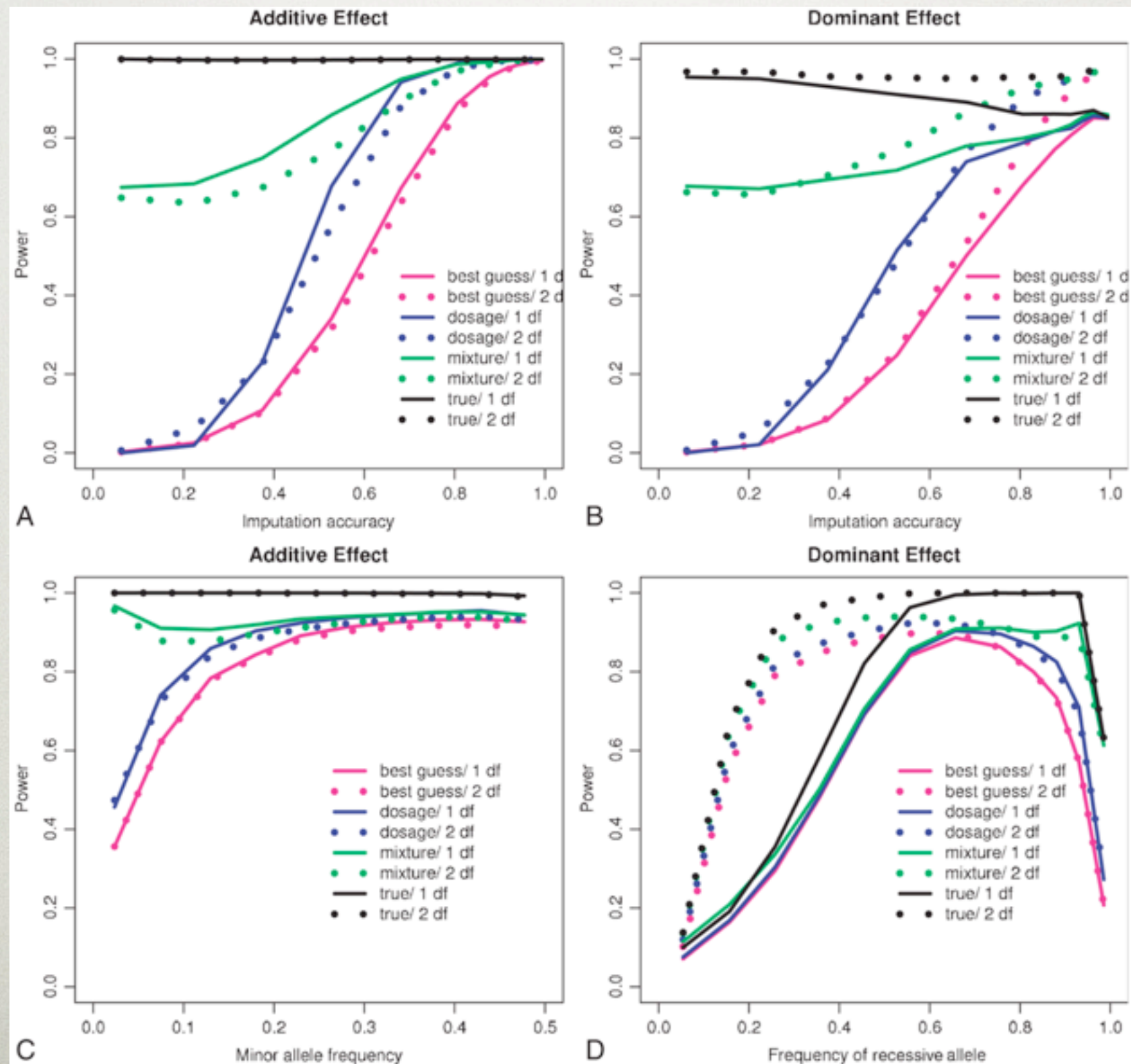


vs. MAF

POWER IN SMALL SAMPLES (LARGE EFFECTS)

vs.
accuracy

vs. MAF



Zheng et al., 2011

SOFTWARE FOR IMPUTATIONS

- MACH / mimimac
- BEAGLE
- IMPUTE2
- BIMBAM
- ...

CONCLUSIONS - ANALYSIS OF IMPUTED DATA

- Using regression onto genotype probabilities is a valid and powerful method for standard scenarios
- Use of ML / mixture method can give extra power in case of small samples and large effects