

A general approach to single-nucleotide polymorphism discovery

G Matrh et al.
1999

SNPs

- a one base DNA sequence variation between two individuals of a same species
- it is the most abundant sequence variation in populations
- almost all SNPs have only two alleles (forms)
- minor allele frequency

data in 1999

- ~1M bp of finished human reference sequence in 10 regions
- data base of express sequence tags(ESTs)
- ESTs are reverse transcribed RNA sequences form different individuals
- length= \sim 300bp

alignment

- first, repeats in the reference is covered
- ESTs are aligned to sequence according to a common anchor
- then, error, gaps, inserts are propagated in the reminding EST
- 1365 hits were found in 147 clusters
- representing 80,469 bp of sequence 38% single, 81% by 8 for fewer ESTs coverage

paralogue

- paralogues are regions highly similar DNA sequences of an individual
- they may have arisen from the same evolutionary root
- must not confuse difference in paralogues with SNPs

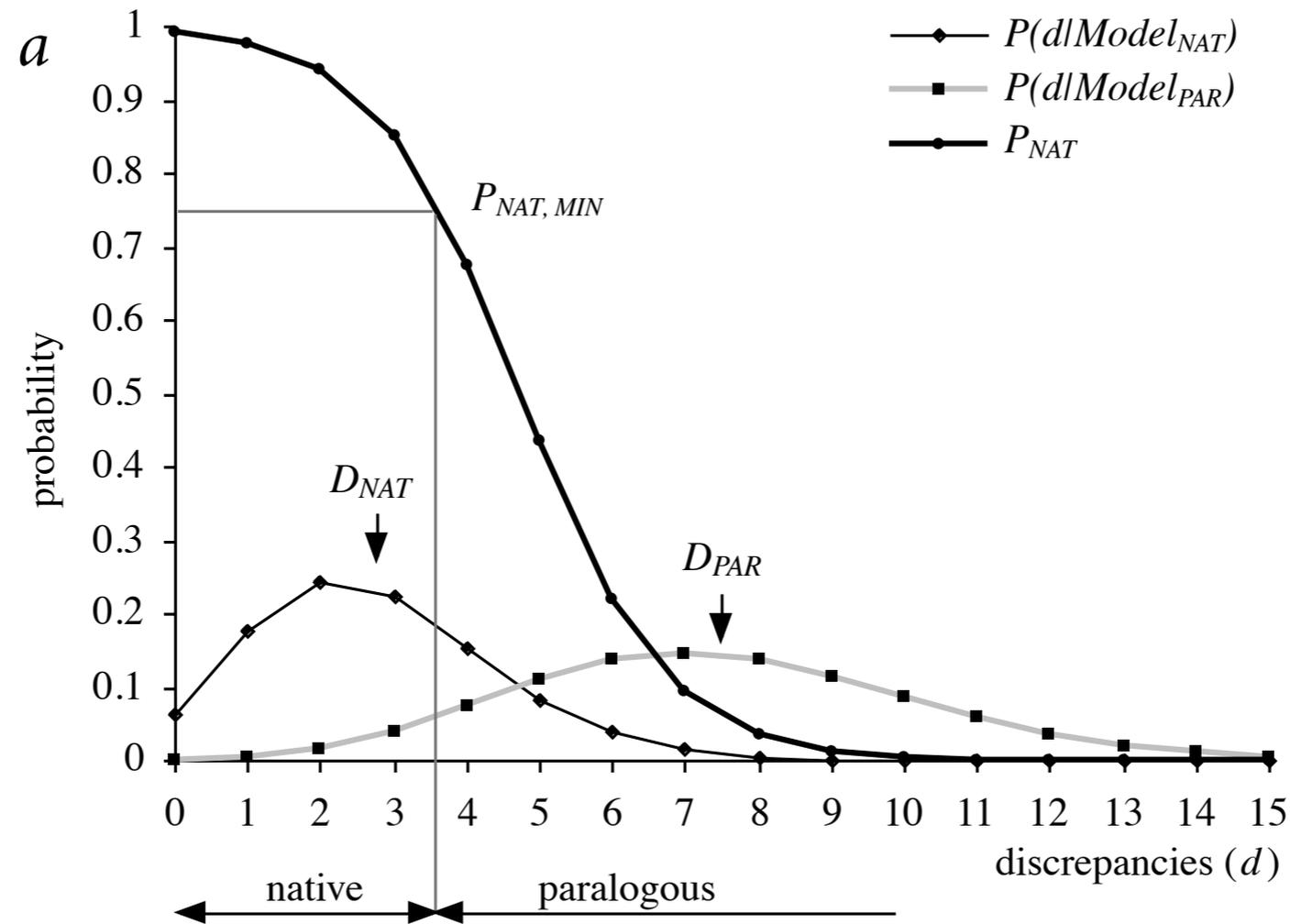
paralogue discrimination

- paralogous sequences have pairwise dissimilarity rate higher than $P_{PAR}=2\%$
- SNP rate is $P_{SNP}=0.1\%$
- they can be differentiated by percent different bases

model for paralogue discrimination

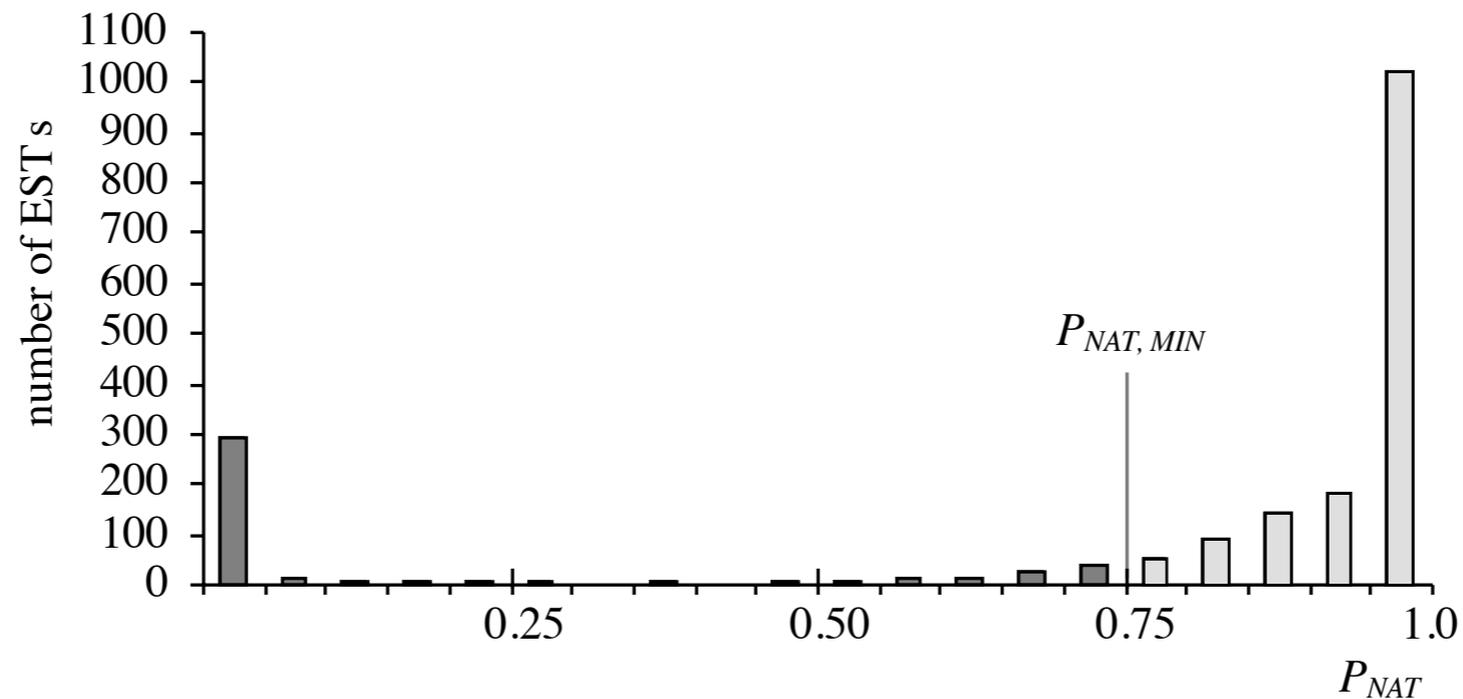
- in a sequence with length L
- expected number of base difference:
from paralogue: $D_{PAR} = L \times P_{PAR} + E$
from SNP: $D_{NAT} = L \times P_{POLY,2} + E$
- They can be approximated as a Poisson distribution
- ratio:
$$P(Model_{NAT} | d) = \frac{1}{1 + e^{(D_{NAT} - D_{PAR})} \cdot \binom{D_{PAR}}{D_{PAR}}}$$

paralogue



hypothetical uniform gene distribution
250 bp long EST $E=2.525$
no mistake in reference

paralogue



experimental results for 1954 cluster
members
(anchored to 10 genomic sequences)

Bayesian SNP inference

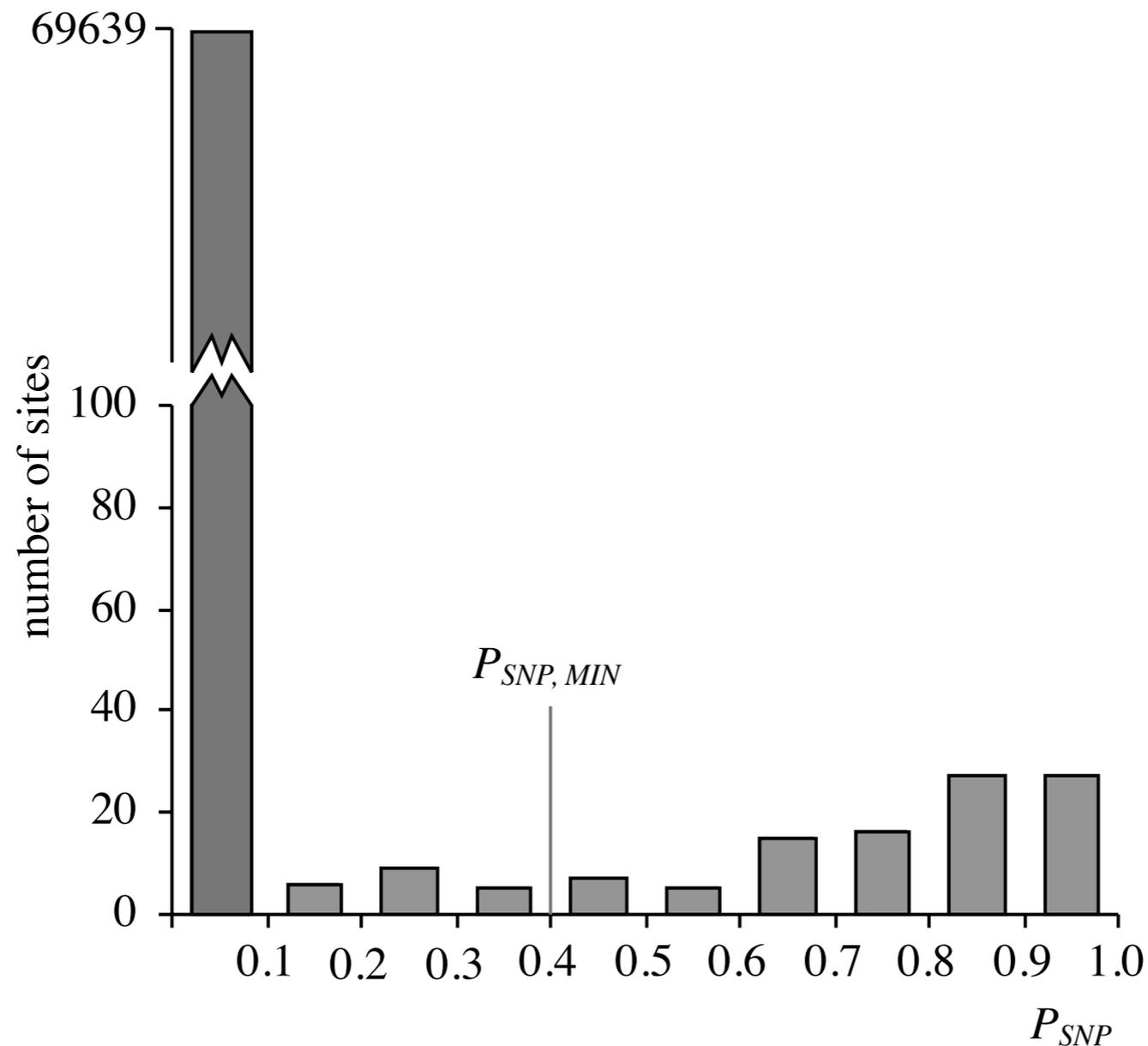
- prior probability for SNP is 0.003
- likelihood of data given sequence (probability of correct sequencing) is estimated from chromatograms
- with these two piece of information, posterior probability is calculated

Bayesian SNP inference

$$P(S_1, \dots, S_N | R_1, \dots, R_N) = \frac{\frac{P(S_1 | R_1)}{P_{\text{Prior}}(S_1)} \cdot \dots \cdot \frac{P(S_N | R_N)}{P_{\text{Prior}}(S_N)} \cdot P_{\text{Prior}}(S_1, \dots, S_N)}{\sum_{\text{every } (S_{i_1}, \dots, S_{i_N})} \frac{P(S_{i_1} | R_1)}{P_{\text{Prior}}(S_{i_1})} \cdot \dots \cdot \frac{P(S_{i_N} | R_N)}{P_{\text{Prior}}(S_{i_N})} \cdot P_{\text{Prior}}(S_{i_1}, \dots, S_{i_N})}$$

- posterior probability for all $S_1 \dots S_N$ are calculated
- SNP probability is the sum of all probability associated with $S_1 \dots S_N$ that have a SNP

Results

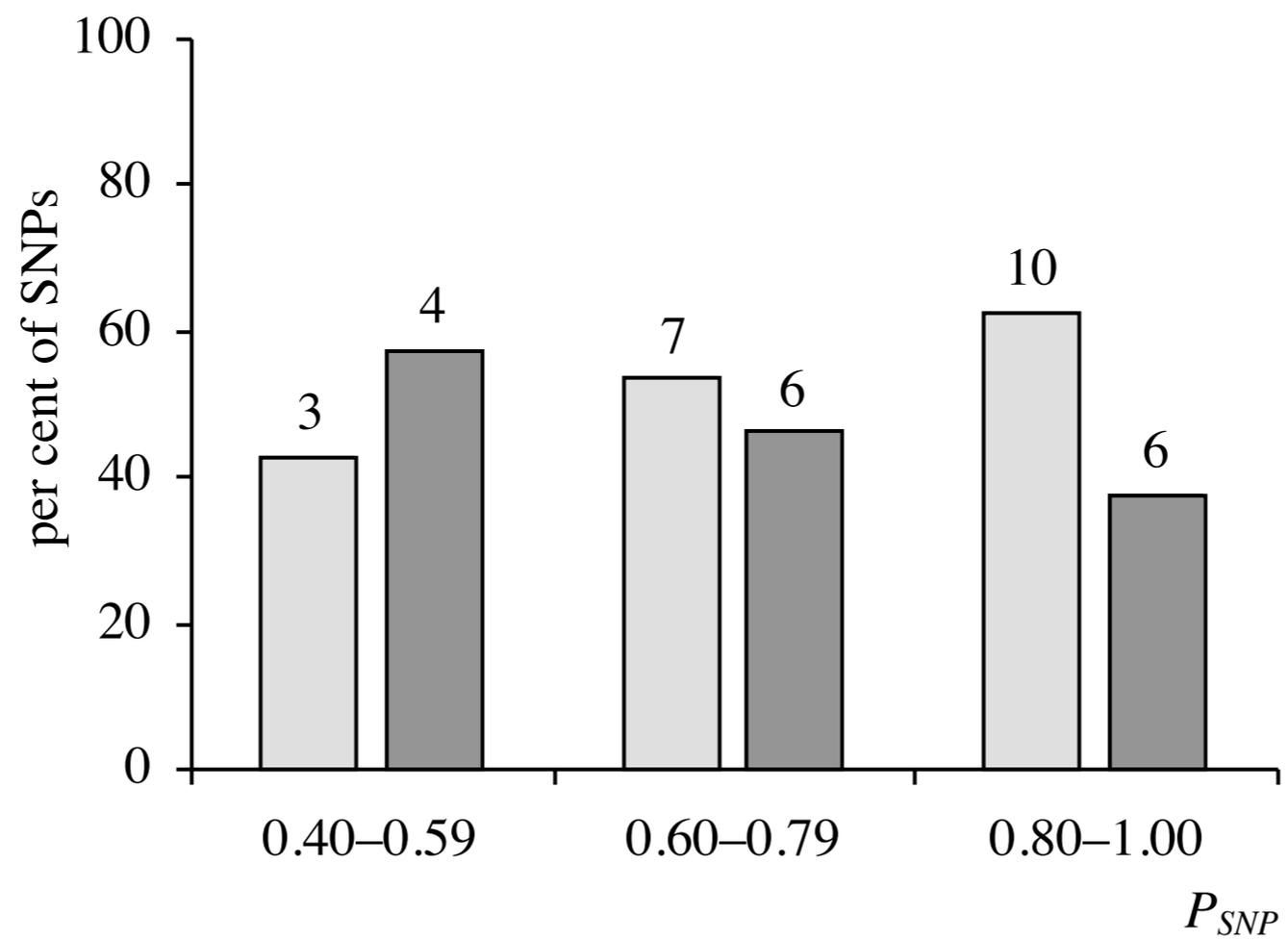


- 97 candidates were found out of 80,469 bp covered sequences

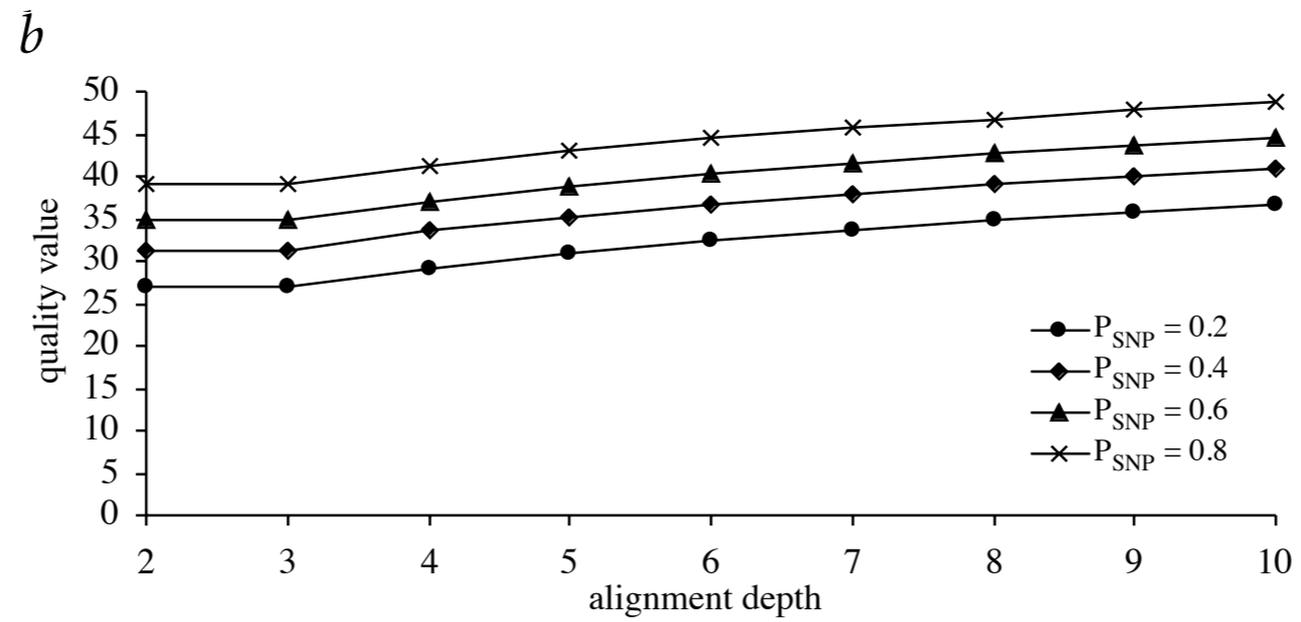
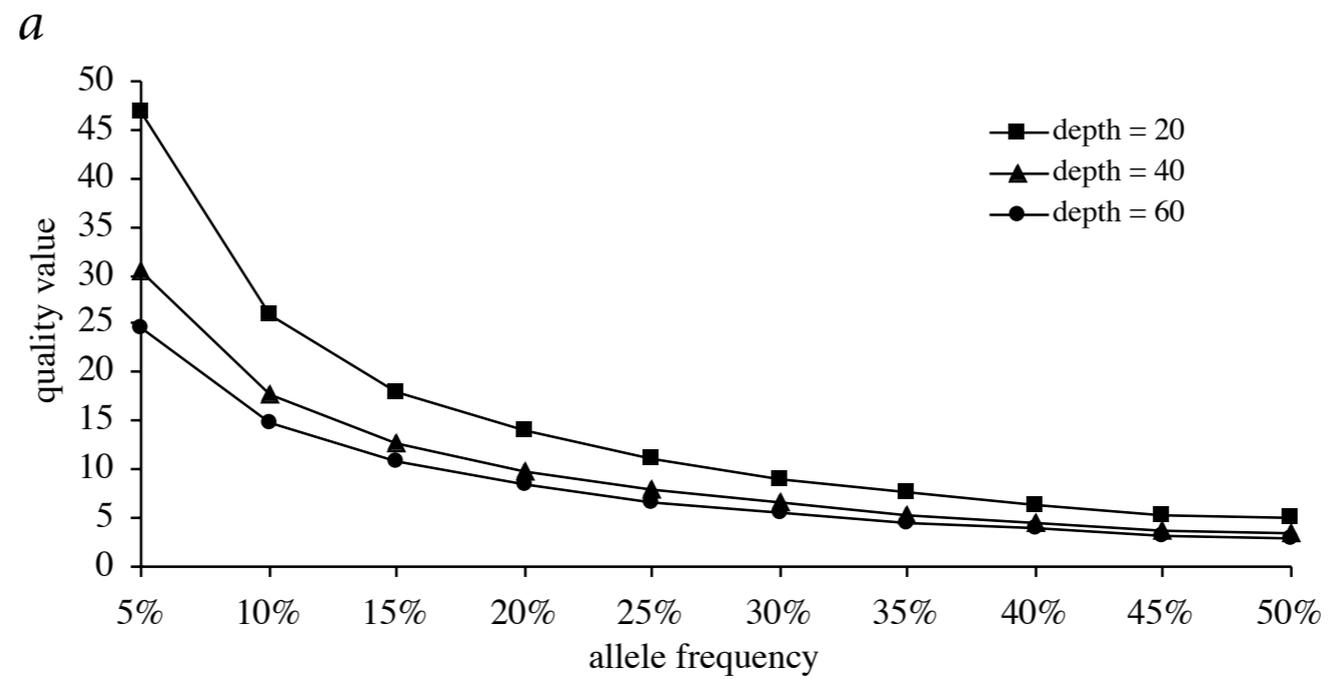
verification

- 38 of these are near the 3' end of ESTs, and verified to have problems with cDNA library construction
- 18 of these can not be analyzed for various reasons
- for the remaining 36 sites, they confirmed 20 sites at least 1 in 4 individuals screened
- overall confirmation rate is reported to be 56%
- another SNP was found in 11,455bp STS

verification



results



conclusion

- Poly-bayes offers a relative straight forward way of finding SNP sites
- reasonable sensitivity and accuracy
- designed for long reads