• Developed in 1990 by Andy Clark

- Developed in 1990 by Andy Clark
- The phasing problem now is related to PCR ambiguity in diploid individuals

- Developed in 1990 by Andy Clark
- The phasing problem now is related to PCR ambiguity in diploid individuals
- For example, suppose that the two chromosomes are:

• How can we resolve this ambiguity?

• If there are k ambiguous sites, there are an exponential (in k) number of possible explanations of the ambiguity

- If there are k ambiguous sites, there are an exponential (in k) number of possible explanations of the ambiguity
- Assume that all sites have *at most* two alleles (the infinite sites model)

- If there are k ambiguous sites, there are an exponential (in k) number of possible explanations of the ambiguity
- Assume that all sites have *at most* two alleles (the infinite sites model)
- Denote the two alleles present by 0 and 1.

```
Haplotypes : 1) 001010

\longrightarrow Genotype : 0(\frac{0}{1})10(\frac{0}{1})0

2) 011000
```

- If there are k ambiguous sites, there are an exponential (in k) number of possible explanations of the ambiguity
- Assume that all sites have *at most* two alleles (the infinite sites model)
- Denote the two alleles present by 0 and 1.

$$\begin{array}{rcl} \textit{Haplotypes}: & 1) & 001010 \\ & & & \longrightarrow & \textit{Genotype}: 0(\frac{0}{1})10(\frac{0}{1})0 \\ & & 2) & 011000 \end{array}$$

• If we were given just the genotype, then there are two possible explanations (k = 2)

0 <u>0</u> 10 <u>0</u> 0	0 <u>0</u> 10 <u>1</u> 0
0 <u>1</u> 10 <u>1</u> 0	0 <u>1</u> 10 <u>0</u> 0

• Notation: A genotype will be a sequence of 0, 1, and 2's [ex. 012120]

- Notation: A genotype will be a sequence of 0, 1, and 2's [ex. 012120]
- '2' denotes ambiguous site

- Notation: A genotype will be a sequence of 0, 1, and 2's [ex. 012120]
- '2' denotes ambiguous site
- haplotype will be a sequences of 0 and 1's [ex. 0101010]

- Notation: A genotype will be a sequence of 0, 1, and 2's [ex. 012120]
- '2' denotes ambiguous site
- haplotype will be a sequences of 0 and 1's [ex. 0101010]
- An explanation of a genotype will be a *pair* of haplotypes

Haplotypes : 1)001010
$$\longrightarrow$$
Genotype : $0(\frac{0}{1})10(\frac{0}{1})0$ 2)0110000(2)10(2)0

- Two easy cases:
 - Homozygous: genotype with only 0's and 1's. The mother and the father chromosome have the same composition.
 - Single heterozygote: genotype with a single '2'. The explanation will still be unique.

- Two easy cases:
 - Homozygous: genotype with only 0's and 1's. The mother and the father chromosome have the same composition.
 - Single heterozygote: genotype with a single '2'. The explanation will still be unique.
- Multiple heterozygote case: More than single '2's in the genotype [ex. 01212]

- Two easy cases:
 - Homozygous: genotype with only 0's and 1's. The mother and the father chromosome have the same composition.
 - Single heterozygote: genotype with a single '2'. The explanation will still be unique.
- Multiple heterozygote case: More than single '2's in the genotype [ex. 01212]
- Number of explanations for a genotype with k ambiguous sites is 2^{k-1} [in this case 01010, 01011, 01110, 01111]

• Suppose that we are given 5 sites on 5 individuals

- Suppose that we are given 5 sites on 5 individuals
- Start with the easy cases, the homozygotes and single heterosygotes

- Suppose that we are given 5 sites on 5 individuals
- Start with the easy cases, the homozygotes and single heterosygotes

Example. Individuals: 01202, 20000, 12121, 00100, and 01000

- Suppose that we are given 5 sites on 5 individuals
- Start with the easy cases, the homozygotes and single heterosygotes

Example. Individuals: 01202, 20000, 12121, 00100, and 01000

• Starting with the 'easy' cases, we are able to resolve 3 of the 5 individuals:

- Suppose that we are given 5 sites on 5 individuals
- Start with the easy cases, the homozygotes and single heterosygotes

Example. Individuals: 01202, 20000, 12121, 00100, and 01000

• Starting with the 'easy' cases, we are able to resolve 3 of the 5 individuals:

Individual	Inferred Mate – Pair
20000 :	10000 - 00000
00100:	00100 - 00100
01000 :	01000 - 01000

• For the unresolved chromosomes (01202, 12121), we have two possible explanations for each:

Individual	Explanation 1	Explanation 2
01202 :	01101	01001
	01000	01100
12121 :	10101	11101
	11111	10111

Now, if we look at the possible explanations together, we see that there is overlap between the inferred mate-pairs from the simple cases and the possible explanations for the ambiguous cases.

Individual	Inferred Mate – Pair
20000 :	10000 - 00000
00100 :	00100 - 00100
01000 :	01000 - 01000

Individual	Explanation 1	Explanation 2
01202 :	01101	01001
	01000	01100
12121 :	10101	11101
	11111	10111

- Notation: ${\sf G}={\sf H}\oplus{\sf H}'$ denotes that haplotype pair (H,H') is an explanation for ${\sf G}$
- Notation: 'H → ^C G' denotes G can be resolved using H ∈ C, i.e., there exists H' s.t. G = H ⊕ H', and we call H' the *inferred haplotype*
- Example: For G=02112:

 $\begin{array}{rl} & \textit{Possible Mate} - \textit{Pairs} \\ G: & 00110 \oplus 01111 \\ & 01110 \oplus 00111 \end{array}$

Clark Algorithm

- Find all homozygotes and single heterozygotes and make a list of all of the haplotypes involved in the unique explanations.
- While there are remaining genotypes that are unresolved, attempt to find a haplotype from the list that helps resolve some unresolved genotype. If such a haplotype exists, add the corresponding mate-pair to the list of haplotypes and label the genotype as resolved.

Clark Algorithm

- Find all homozygotes and single heterozygotes and make a list of all of the haplotypes involved in the unique explanations.
- While there are remaining genotypes that are unresolved, attempt to find a haplotype from the list that helps resolve some unresolved genotype. If such a haplotype exists, add the corresponding mate-pair to the list of haplotypes and label the genotype as resolved.

Example: If the genotype is 12121, then 10101 is the *mate-pair* of 11111 in the explanation of genotype 12121

• If we are able to find a haplotype in the list and thus explain the ambiguous haplotype, we have inferred a uniquely inferred additional haplotype, that of the mate-pair.

• Drawbacks to the Clark Methods:

- Drawbacks to the Clark Methods:
 - Possible that a genotype can be resolved in multiple ways from the list, thus yielding multiple explanations. Which one do you choose? (Associated problem: anomalous genotypes i.e., explanations that are not correct)

- Drawbacks to the Clark Methods:
 - Possible that a genotype can be resolved in multiple ways from the list, thus yielding multiple explanations. Which one do you choose? (Associated problem: anomalous genotypes i.e., explanations that are not correct)
 - Possible that unresolved genotypes are all incompatible with haplotypes in list. Then the algorithm will stop and leave behind "orphan", or unresolved genotypes

- Drawbacks to the Clark Methods:
 - Possible that a genotype can be resolved in multiple ways from the list, thus yielding multiple explanations. Which one do you choose? (Associated problem: anomalous genotypes i.e., explanations that are not correct)
 - Possible that unresolved genotypes are all incompatible with haplotypes in list. Then the algorithm will stop and leave behind "orphan", or unresolved genotypes
 - If there are no homozygotes or single heterozygotes, then the algorithm cannot start.

- Drawbacks to the Clark Methods:
 - Possible that a genotype can be resolved in multiple ways from the list, thus yielding multiple explanations. Which one do you choose? (Associated problem: anomalous genotypes i.e., explanations that are not correct)
 - Possible that unresolved genotypes are all incompatible with haplotypes in list. Then the algorithm will stop and leave behind "orphan", or unresolved genotypes
 - If there are no homozygotes or single heterozygotes, then the algorithm cannot start.
 - The order in which you resolve the haplotypes matters. A different ordering may produce a different haplotype list.

- Drawbacks to the Clark Methods:
 - Possible that a genotype can be resolved in multiple ways from the list, thus yielding multiple explanations. Which one do you choose? (Associated problem: anomalous genotypes i.e., explanations that are not correct)
 - Possible that unresolved genotypes are all incompatible with haplotypes in list. Then the algorithm will stop and leave behind "orphan", or unresolved genotypes
 - If there are no homozygotes or single heterozygotes, then the algorithm cannot start.
 - The order in which you resolve the haplotypes matters. A different ordering may produce a different haplotype list.
- What is the probability of the algorithm stopping prematurely or not being able to start?

- How can we estimate the number of 2's?
- Population model: Infinite Sites Model (at most *one* mutation can happen at any site on the chromosome)

- How can we estimate the number of 2's?
- Population model: Infinite Sites Model (at most *one* mutation can happen at any site on the chromosome)
- Neutral Model of Evolution [

- How can we estimate the number of 2's?
- Population model: Infinite Sites Model (at most *one* mutation can happen at any site on the chromosome)
- Neutral Model of Evolution [
- expected number of mismatches of a DNA sequence:

$$\Theta = L\Theta_{nt}$$

where

• L =length of the sequence

•
$$\Theta_{nt} = 4N\mu$$

- N = the effective population size
- μ is the mutation rate per nucleotide per generation

• For example, in Drosophila, $\Theta = .005L$.

$$Pr(2 \text{ sequences have m mismatches}) = rac{1}{ heta+1} \left(rac{ heta}{ heta+1}
ight)^m$$

• In the infinite sites model,

$$Prob(2 ext{ genes identical}) = rac{1}{1+ heta}(m=0)$$

• Probability of two different genes is

$$1 - {\it Prob}(2 \; {
m genes} \; {
m identical}) = { heta \over 1 + heta}$$

• Now, if we have N diploid individuals,

• Now, if we have N diploid individuals,

$$Pr(No \ Homozygotes) = rac{ heta^3 + 4 heta^2 + 2 heta}{(1+ heta)(2+ heta)(3+ heta)}$$

• Now, if we have N diploid individuals,

$$Pr(No \ Homozygotes) = rac{ heta^3 + 4 heta^2 + 2 heta}{(1+ heta)(2+ heta)(3+ heta)}$$
 $Pr(Single \ site \ Heterozygote) = rac{ heta}{(1+ heta)^2}$

• Now, if we have N diploid individuals,

$$Pr(No \ Homozygotes) = \frac{\theta^3 + 4\theta^2 + 2\theta}{(1+\theta)(2+\theta)(3+\theta)}$$
$$Pr(Single \ site \ Heterozygote) = \frac{\theta}{(1+\theta)^2}$$
$$Pr(Algorithm \ won't \ start) = \left[1 - \frac{1}{1+\theta} - \frac{\theta}{(1+\theta)^2}\right]^N$$

• Now, if we have N diploid individuals,

Pr(No homozygotes) [i.e., Pr(algorithm won't start)] is obtained using *Ewing's Sampling Lemma*

$$Pr(No \ Homozygotes) = \frac{\theta^{3} + 4\theta^{2} + 2\theta}{(1+\theta)(2+\theta)(3+\theta)}$$
$$Pr(Single \ site \ Heterozygote) = \frac{\theta}{(1+\theta)^{2}}$$
$$Pr(Algorithm \ won't \ start) = \left[1 - \frac{1}{1+\theta} - \frac{\theta}{(1+\theta)^{2}}\right]^{N}$$

Thus, if $\Theta > 0.5$, Clark's algorithm will work well.

- How can we deal with the 'orphan' genotypes? What is the probability of not finishing?
- The answer depends on the algorithm implementation and the order in which you resolve the haplotypes.

- How can we deal with the 'orphan' genotypes? What is the probability of not finishing?
- The answer depends on the algorithm implementation and the order in which you resolve the haplotypes.
- Maximum Resolution Problem:

Input: A set of vectors, ambiguous (0,1,2) and resolved (0,1) [genotypes] **Output:** Maximum number of ambiguous vectors that can be

- resolved by successive applications of Clark's Rule.
- Equivalent to minimize the number of orphan genotypes

Theorem

The Maximum Resolution Problem is NP-complete.

Theorem

The Maximum Resolution Problem is NP-complete.

The proof for this theorem is based on a reduction to the satisfiability problem (SAT).

Proof.

- Let X₁, X₂, X₃, and X₄ be variables, and let (x₁, x₁, x₂, x₂, x₃, x₃, x₄, x₄) be literals. The satisfiability problem is
 - Input: Boolean function F
 - **Output**: Is there a truth assignment that makes F true? If so, find such an assignment.

Proof.

- Let X₁, X₂, X₃, and X₄ be variables, and let (x₁, x₁, x₂, x₂, x₃, x₃, x₄, x₄) be literals. The satisfiability problem is
 - Input: Boolean function F
 - **Output**: Is there a truth assignment that makes F true? If so, find such an assignment.
- We are going to attempt to create a 1-1 correspondence between the boolean logic and haplotypes using the above function F.

 To the function F, we associate a set of genotypes (corresponding to columns in the matrix). The number of rows in the matrix will be # variables + 2 × #variables +#clauses +1 = 3V +C +1.

- To the function F, we associate a set of genotypes (corresponding to columns in the matrix). The number of rows in the matrix will be # variables + 2 × #variables +#clauses +1 = 3V +C +1.
- We fill in rows C₁, C₂, C₃ in the following manner: if X_i is absent from a clause, put '1' in both T_i and F_i. If X_i is present in a clause, place '1' in T_i and '0' in F_i. If X_i is present in a clause, then place a '0' in T_i and a '1' in F_i.

- To the function F, we associate a set of genotypes (corresponding to columns in the matrix). The number of rows in the matrix will be # variables + 2 × #variables +#clauses +1 = 3V +C +1.
- We fill in rows C_1 , C_2 , C_3 in the following manner: if X_i is absent from a clause, put '1' in both T_i and F_i . If X_i is present in a clause, place '1' in T_i and '0' in F_i . If $\overline{X_i}$ is present in a clause, then place a '0' in T_i and a '1' in F_i .

• Fill in columns
$$S_1 \rightarrow S_4$$
 by:

• To fill in the $C_1 \rightarrow C_3$ columns, we look at our function F.

-∢ ≣ ▶

- To fill in the $C_1 \rightarrow C_3$ columns, we look at our function F.
- We see that x_1 appears in both clause 2 and clause 3. So we place a '2' in both C_2 and C_3 .

- To fill in the $C_1 \rightarrow C_3$ columns, we look at our function F.
- We see that x₁ appears in both clause 2 and clause 3. So we place a '2' in both C₂ and C₃.
- For the rows containing the literals $(x_1\bar{x_1}, \text{ etc.}) x_1$ appears in C_2 but no other clause.

- To fill in the $C_1 \rightarrow C_3$ columns, we look at our function F.
- We see that x₁ appears in both clause 2 and clause 3. So we place a '2' in both C₂ and C₃.
- For the rows containing the literals $(x_1\bar{x_1}, \text{ etc.}) x_1$ appears in C_2 but no other clause.
- Place a '2' in C_2 . $\bar{x_1}$ appears in C_3 , but no other clause. Place a '2' in C_3 . etc. (Note: the rows x_1 and $\bar{x_1}$ should sum to the X_1 row in the C columns.)

How do we fill in the bottom right corner of the table?

- We have a clause set $C_1... C_c$, one for each clause. All will be ambiguous vectors
- For each K=1,2,...c the first v positions of the vector C_k are zero except for position i such that either x_i or \bar{x}_i appears in C_k . We are blind to the actual truth value
- For the next 2v positions [the literals], place a zero except for any position v+2i-1 where x_i appears in clause C_k or position v+2i where \bar{x}_i appears in clause C_k. These positions are set to 2
- For each r from 1 to c, position 3v+r of $C_k = 0$ if and only if r=k [the diagonal].
- For $r \neq k$ position 3v+r = 2 if and only if clause k and r contain a variable in common [not necessarily a literal in common].
- Otherwise, position 3v+r = 1. [This assignment captures the ambiguity related to the literals that are contained in multiple clauses. You may have both x_i and \bar{x}_i present in F].

We reduce SAT to MR. Start with a generic boolean formula F with C clauses and V variables $X_1 \dots X_v$. F takes a set of vectors V(F) that are ambiguous and resolved as input to the MR problem [the columns of out table drawn above]. We want to show that F has a satisfying truth assignment *iff* V(F) has the maximum number of ambiguous vectors explained by a series of Clark Rules. Recall that $F = (x_2 \lor x_3 \lor \bar{x_4}) \land (x_1 \lor \bar{x_2} \lor x_4) \land (\bar{x_1} \lor \bar{x_3})$. How can we interpret resolution by the Clark Method?

• Suppose that you pick colum T_1 to resolve one of the other columns; this is interpreted as setting the literal X_1 to true.

How can we interpret resolution by the Clark Method?

- Suppose that you pick colum T_1 to resolve one of the other columns; this is interpreted as setting the literal X_1 to true.
- Suppose that we try to resolve column [haplotype] S_1 . S_1 can be resolved by using either the T_1 or F_1 columns. By picking one of the columns, we are fixing X_1 as true or false.
- Using the haplotypes obtained in the resolution of the S columns along with the T and F columns, can we resolve the columns $C_1 \rightarrow C_3$? (NO)

	T_1	F_1	T_2	F_2	<i>T</i> ₃	F ₃	T_4	F_4	<i>S</i> ₁	S_2	S_3	S_4	<i>C</i> ₁	C_2	C3
X_1	1	1	0	0	0	0	0	0	1	0	0	0	0	2	2
X2	0	0	1	1	0	0	0	0	0	1	0	0	2	2	0
X3	0	0	0	0	1	1	0	0	0	0	1	0	2	0	2
X_4	0	0	0	0	0	0	1	1	0	0	0	1	2	2	0
x ₁	0	0	0	0	0	0	0	0	0	0	0	0	0	2	0
xī1	0	0	0	0	0	0	0	0	0	0	0	0	0	0	2
X2	0	0	0	0	0	0	0	0	0	0	0	0	2	0	0
x2	0	0	0	0	0	0	0	0	0	0	0	0	0	2	0
x3	0	0	0	0	0	0	0	0	0	0	0	0	2	0	0
x3	0	0	0	0	0	0	0	0	0	0	0	0	0	0	2
X4	0	0	0	0	0	0	0	0	0	0	0	0	0	2	0
x4	0	0	0	0	0	0	0	0	0	0	0	0	2	0	0
C ₁	1	1	1	0	1	0	0	1	1	2	2	2	0	2	2
C ₂	1	0	0	1	1	1	1	0	2	2	1	2	2	0	2
C3	0	1	1	1	0	1	1	1	2	1	2	1	2	2	0
C_B	0	0	0	0	0	0	0	0	2	2	2	2	1	1	1

Input: A set of unresolved genotypes **Output:** Maximum number of ambiguous vectors that can be resolved by successive applications of Clark's Rule. Input: A set of unresolved genotypes

Output: Maximum number of ambiguous vectors that can be resolved by successive applications of Clark's Rule.

• Recall the notation $R \rightarrow^{C} A$ where R = resolved haplotype and A = ambiguous haplotype.

Input: A set of unresolved genotypes

Output: Maximum number of ambiguous vectors that can be resolved by successive applications of Clark's Rule.

- Recall the notation $R \rightarrow^{C} A$ where R = resolved haplotype and A = ambiguous haplotype.
- If we call A[i] the i-th site of A and R[j] the j-th site in R then the notation means that if A[i] = 0 or 1 then R[i] = 0 or 1 and R[i] = A[i].

Further explanation of the NP-completeness proof: For a more full presentation, see Dan Gusfield's paper

Inference of Haplotypes from Samples of Diploid Populations: Complexity and Algorithms, J. Computational Biology August 2001.

Remarks (Table Setup)

- For every variable, there are two columns T and F
- A column of selectors exists for each random variable
- A column exists for each clause
- The first V rows are associated with the random variables
- The next set of rows are associated with the literals
- The final set of rows corresponds to clauses, including the mysterious (for a little while longer anyways) C_b
- $T_1, F_1, ..., T_4, F_4$ are all *resolved* columns while the rest are *unresolved*

Remarks (General Properties)

- $T_i \rightarrow^C S_i$ or $F_i \rightarrow^C S_i$ but $T_i \not\rightarrow^C S_j$ and $F_i \not\rightarrow^C S_j$. i.e. Column T_i can be applied by Clark Rule to column S_i but to no other selector column.
- At most one T or F can be applied to any S. i.e. set X₁ to either T or F, but not both!
- We interpret as follows: if $T_i \rightarrow^C S_i$ as ' $X_i = \text{true}$ '; if $F_i \rightarrow^C S_i$ as ' $X_i = \text{false}$ '. Suppose $T_1 \rightarrow^C S_1 F_2 \rightarrow^C S_2$ $F_3 \rightarrow^C S_3 T_4 \rightarrow^C S_4$, so the *inferred* vectors are $R_1 = S_1 \ominus T_1$, $R_2 = S_2 \ominus F_2$, $R_3 = S_3 \ominus F_3$, $R_4 = S_4 \ominus T_4$
- R₁, R₂, R₃, R₄ can be applied only to C₁, C₂ and C₃. [Consider the Blocking Clause C_b] The last entry of R₁ to R₄ will be a 1
- No T or F can be applied to the C vectors [Because of the Blocking Clause, *C_b*]
- $T_i \rightarrow^C S_i$, then $R_k \rightarrow^C C_k$ iff the literal x_i appears in C_k . Similarly, $F_i \rightarrow^C S_i$ then $R_k \rightarrow^C C_k$ iff the literal \bar{x}_i appears in C_k .