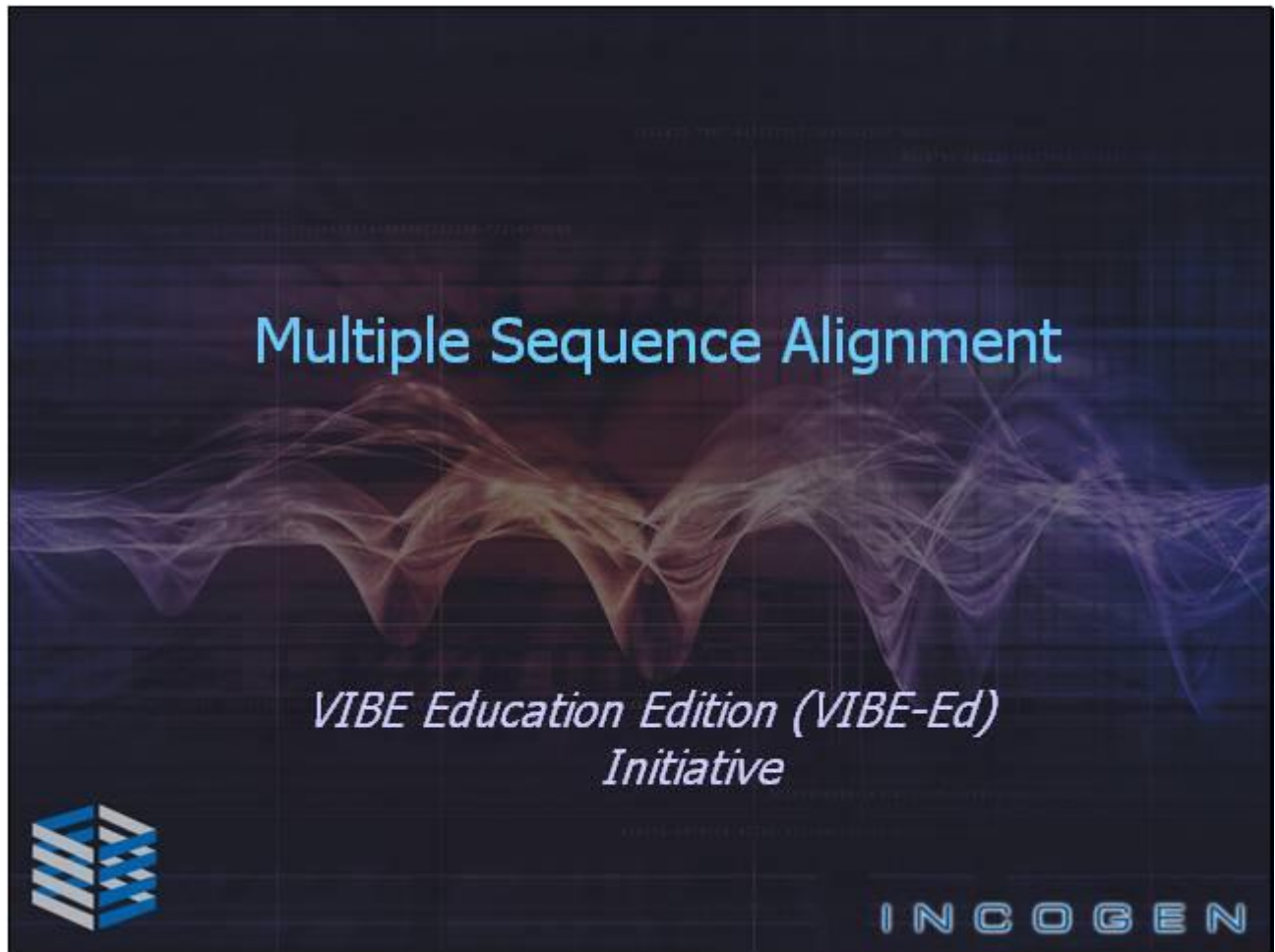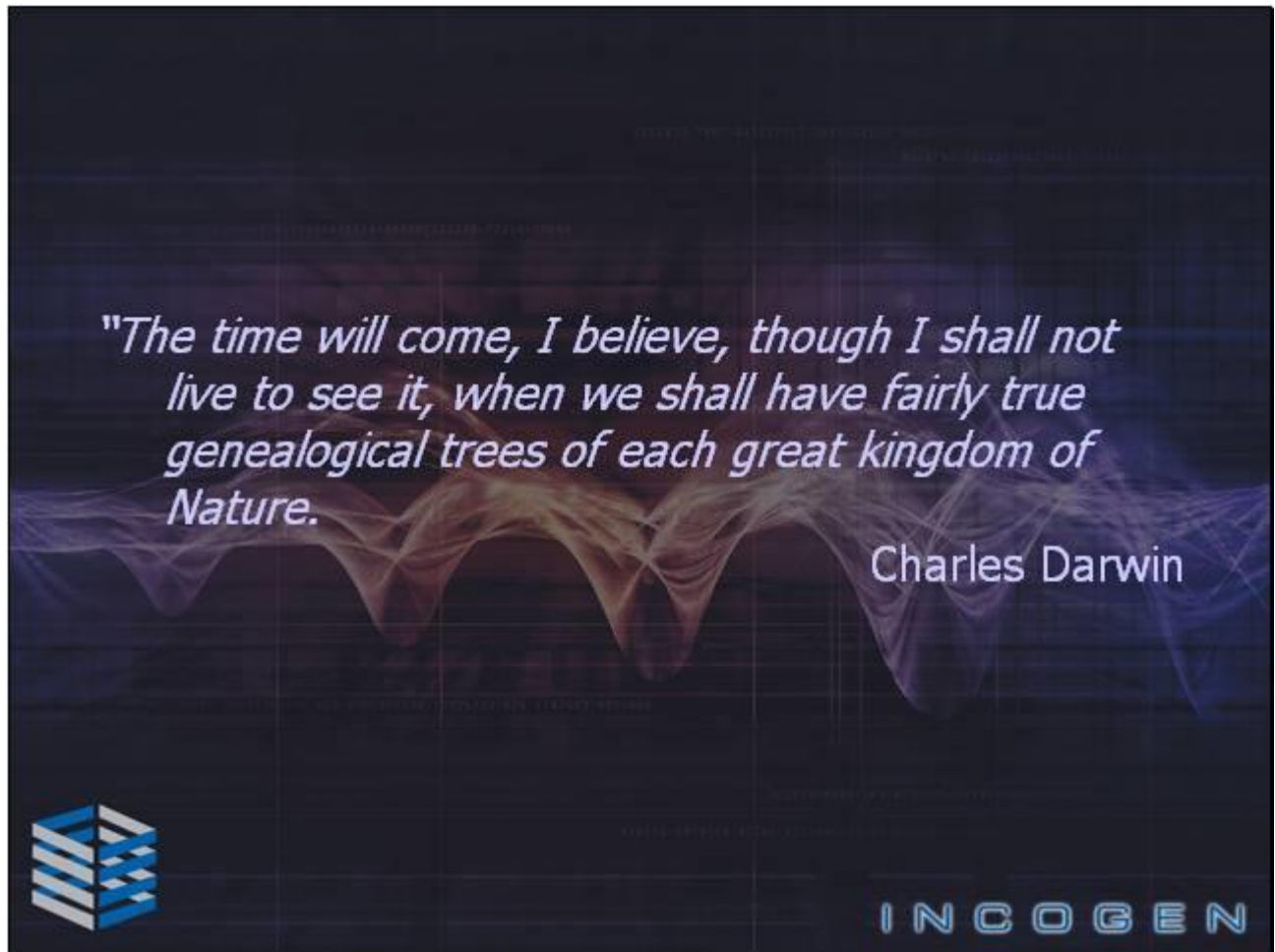**Slide 1 - Multiple Sequence Alignment**



**Slide notes**

This presentation will give you some information about multiple sequence alignments.
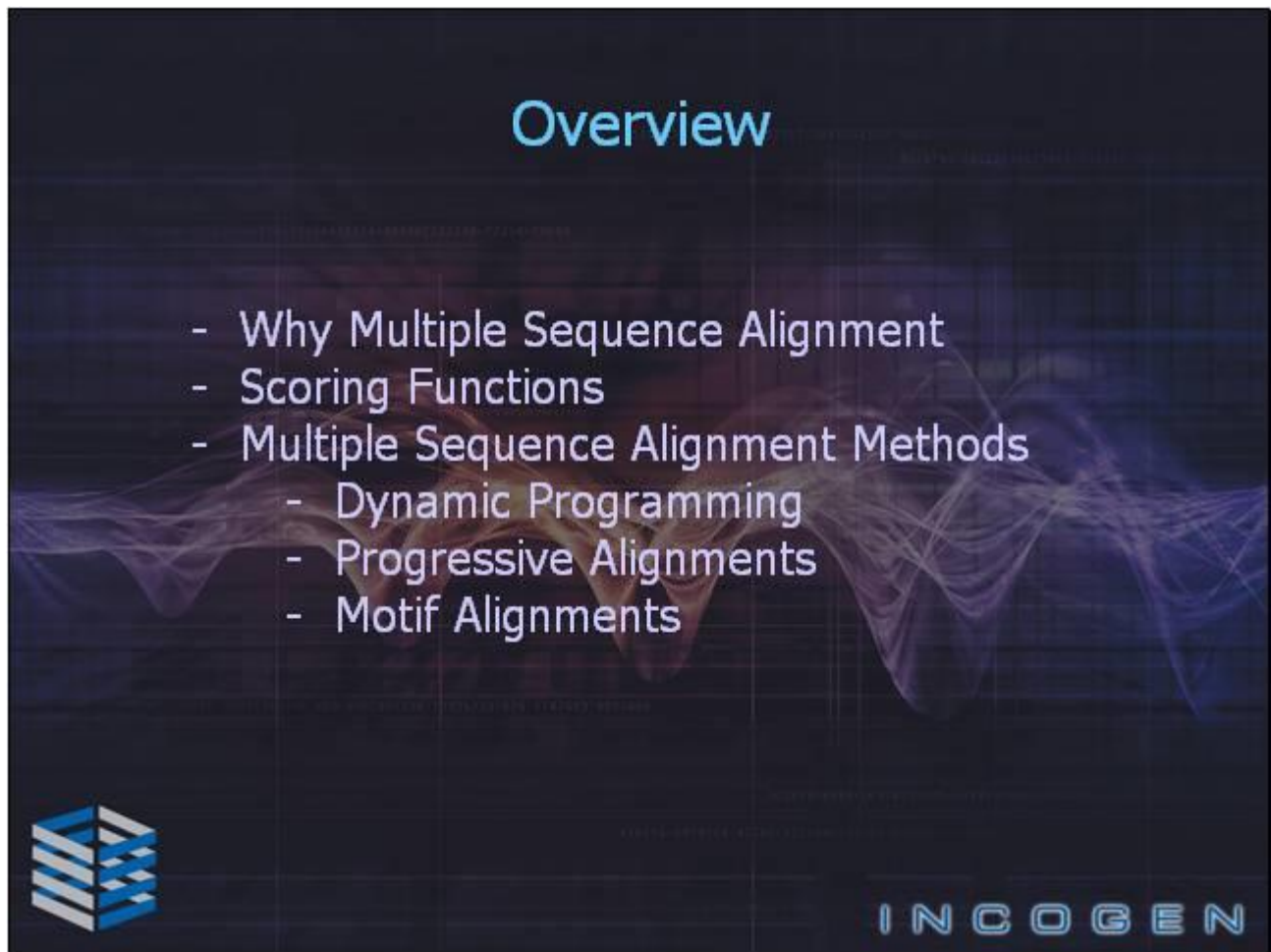
**Slide 2 - Slide 2**



**Slide notes**

Charles Darwin once said, "The time will come, I believe, though I shall not live to see it, when we shall have fairly true genealogical trees of each great kingdom of Nature."

**Slide 3 - Overview**



**Slide notes**

Over the next few minutes, we'll discuss why we need multiple sequence alignments and then how we calculate and score them.

**Slide 4 - Multiple alignment**



**Slide notes**

Why do we need multiple sequence alignments?  What information can we learn from them that we can't learn from pairwise sequence alignments?  Typically with a pairwise alignment, we infer biological relationships from the sequence similarity.  With a multiple alignment, we know that the sequences are biologically related, and we use the multiple alignment to find the areas of sequence similarity that could point to the structure of an evolutionary ancestor or provide information about the evolutionary history of the sequences.  Multiple sequence alignments, or MSAs, are also more sensitive to sequence similarities than a pairwise alignment because the conserved regions could be so dispersed that a pairwise alignment wouldn't find them.

**Slide 5 - Why do we care about multiple sequence alignment?**



**Slide notes**

So, MSAs allow us to infer phylogenetic relationships. They can also help us to elucidate biological facts about proteins since most conserved regions are biologically significant. MSAs can also help us to formulate and test hypotheses about protein 3-D structure and function.
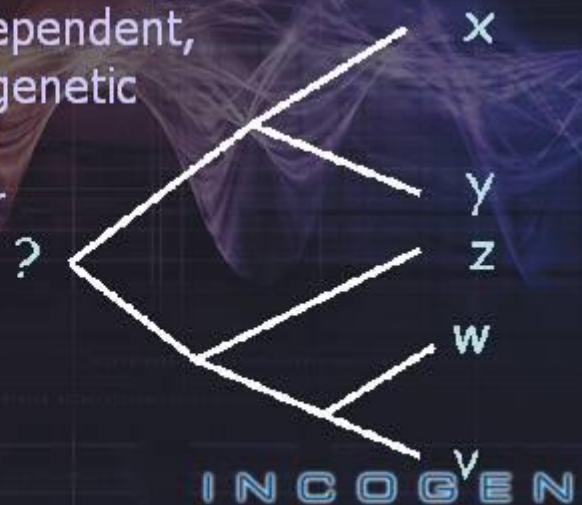
**Slide 6 - Multiple Sequence Alignment (MSA) Defined**



**Slide notes**

An MSA is the alignment of more than two protein or nucleotide sequences. The alignments in an MSA are global, and gaps are added as the sequences are aligned so that all of the sequences have the same length.

**Slide 7 - Scoring Function**



**Slide notes**

In order to score an alignment, we have to be able to quantitatively calculate how good it is.  Any scoring algorithm we use needs to take into account that some positions are more conserved than others, which is called position-specific scoring, and that the sequences are biologically related by a phylogenetic tree.

**Slide 8 - Scoring Functions**



**Slide notes**

All columns in the alignment are treated as statistically independent.

**Slide 9 - Scoring Function:  Definitions**



Scoring Function: Definitions

- Define m

$$m = \begin{matrix} AC-GCGG-C \\ AC-GC-GAG \\ GCACC-GAG \end{matrix}$$

- $m_i^j$ = symbol in column i for sequence j
  - $m_4^2 = G$

- $c_{ia}$ = observed counts for residue a in column i
  - $c_{1A} = 2$, $c_{1C} = 0$, $c_{1G} = 1$, $c_{1T} = 0$, $c_{1-} = 0$

INCOGEN

**Slide notes**

Let's go through some terminology first.  The alignment is referred to as m.  The letter i typically refers to the column number, j refers to the sequence, and a refers to the specific residue. There are different scoring functions that can be used to calculate an MSA.

**Slide 10 - Scoring Function: Minimum Entropy**



## Scoring Function: Minimum Entropy

- Probability of column $m_i$:
$$P(m_i) = \prod_a (p_{ia})^{c_{ia}}$$

- Define column score as:
$$S(m_i) = -\log\left[\,P(m_i)\,\right]$$
$$= -\log\left[\,\prod_a (p_{ia})^{c_{ia}}\,\right]$$
$$= -c_{ia}\log\left[\,\prod_a (p_{ia})\,\right]$$
$$= -\sum_a c_{ia}\log (p_{ia})$$

- Measures variability observed in an aligned columns of residue

- Estimate for $\quad p_{ia} = \dfrac{c_{ia}}{\sum_{a'} c_{ia'}}$

- "Good alignment" minimizes total entropy $\sum_i S(m_i)$

INCOGEN

**Slide notes**

The first scoring function we are going to discuss is the minimum entropy scoring function. The goal of this scoring algorithm is to minimize the entropy, or randomness, in the in the alignment. To calculate the entropy of the alignment, first, we must calculate the probability of column i and then use that probability to calculate a score for that column. This score measures the variability observed in the aligned column i. By minimizing the sum of this column score over all of the columns, we minimize the entropy and create a "good" alignment.

**Slide 11 - Scoring Function: Minimum Entropy Example**



**Slide notes**

Let's look at a simple example. Here is a short multiple alignment consisting of three nucleotide sequences and nine columns.

**Slide 12 - Scoring Function:  Sum Of Pairs**



**Scoring Function: Sum Of Pairs**

- The sum-of-pairs (SP) score of a multiple alignment m is the sum of the scores of all induced pairwise alignments.
- SP score for column $m_i$ is:

$$S(m_i) = \sum_{k<l} s(m_i^k, m_i^l)$$

$s(a,b)$ is obtained from substitution matrix

**Slide notes**

Another scoring algorithm is the sum of pairs algorithm.  In this scoring algorithm, the score of an MSA is the sum of the scores of all of the pairwise alignments.

**Slide 13 - Notation**



**Slide notes**

The notation for this scoring algorithm is fairly straight forward. At first glance, we want to do a sum of the alignment scores between every sequence k and every sequence l. First, we take the sum over every sequence l and then we take the sum over every sequence k. When we do that, however, we find that we are including the alignments of all of the sequences with themselves, such as (1,1), and that we are double-counting all of the other alignments, since the alignment of sequence 1 with sequence 2 is the same as the alignment of sequence 2 with sequence 1.

**Slide 14 - Notation**



**Slide notes**

What we really want is the sum of the alignment scores between every sequence k and every other sequence l, where k < l.  This will prevent counting alignments between a sequence and itself, and it will prevent double-counting the other alignments,

**Slide 15 - Notation**



**Slide notes**

…leaving us with this formula.

**Slide 16 - Scoring Function:  Sum Of Pairs Example**



**Slide notes**

Let's look at an example.  For this particular example, we have used BLOSUM50 as the pairwise sequence alignment scoring matrix.

**Slide 17 - Multiple Alignment Methods**



**Slide notes**

Now that we have a scoring function, let's take a brief look at the methods that use these functions: dynamic programming, heuristic, progressive, progressive with refinement, and model or profile alignment.

**Slide 18 - Dynamic Programming (Optimal Solution)**



**Slide notes**

Let's look at dynamic programming first. Since this is the method used to find pairwise alignments, it seems like an obvious first choice to calculate multiple alignments. The modifications to the pairwise alignment algorithm are fairly straightforward. We now have k sequences, where k > 2, and the dynamic programming array now has N dimensions instead of two. The calculations themselves remain the same.

**Slide 19 - Slide 19**



**Slide notes**

This is an example of what the dynamic programming array looks like when aligning two sequences. The complexity of this algorithm is O(n2), where n is the length of the sequence.

**Slide 20 - Slide 20**



**Slide notes**

If we move to a multiple sequence alignment with three sequences, our dynamic programming array looks like a cube. Visually, this is a bit more confusing to look at, and the complexity increased to O(n3).
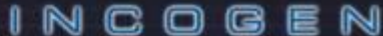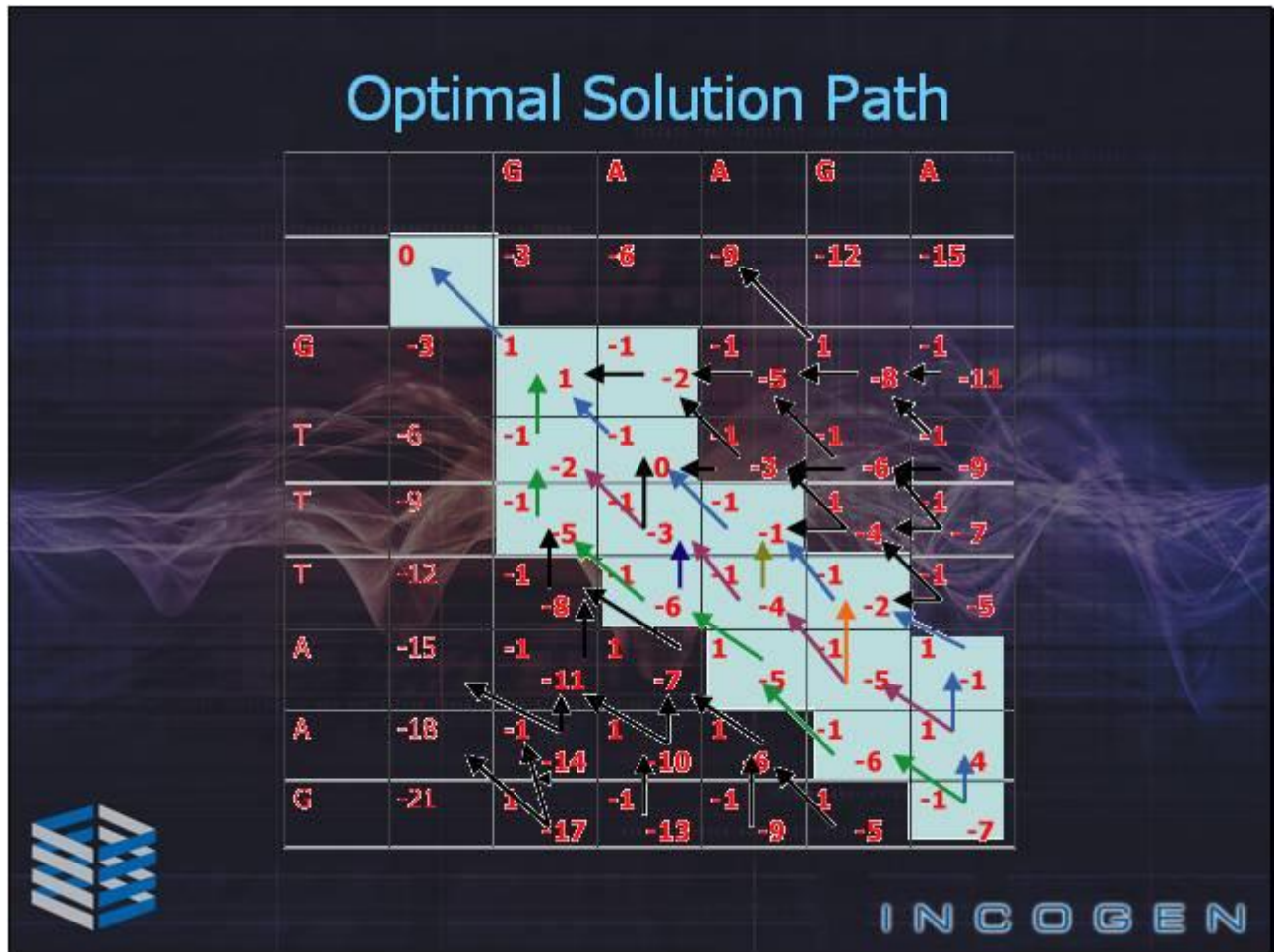
**Slide 21 - Dynamic Programming**



# Dynamic Programming

- Complexity:
  - $O(n^k)$, for k sequences, each n residues long
- Assume sequences of length 300:
  - 2 sequences: 300*300 comparisons ($9*10^4$)
  - 3 sequences: 300*300*300 comparisons ($2.7*10^7$)
  - 4 sequences: $8.1*10^9$
  - 5 sequences: $2.4*10^{12}$
  - 10 sequences: $5.9*10^{24}$
  - 20 sequences: $3.5*10^{49}$
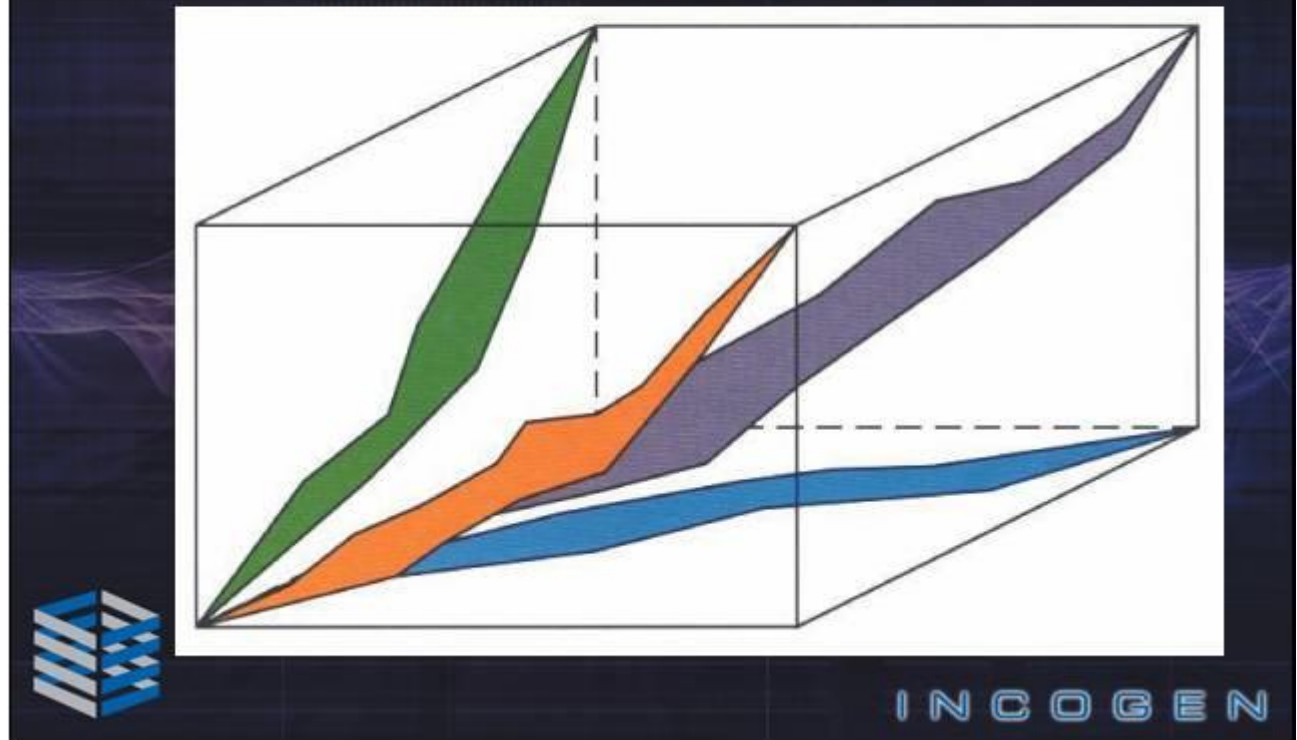  - 30 sequences: $2.1*10^{74}$

INCOGEN

**Slide notes**

So, if we have k sequences, each n residues long, the complexity of the dynamic programming algorithm is O(nk). This grows very quickly ask increases. Let me show you a few examples using sequences of length 300.

**Slide 22 - Slide 22**



**Slide notes**

We can decrease the number of comparisons that we need to do if we remember that for global alignments, the solution is generally found within a small area around the diagonal of the dynamic programming array. We can use a heuristic method to eliminate the areas where solutions are rarely found. This eliminates a lot of generally unneeded calculations at the expense of the rarely found solution that lies outside this region.
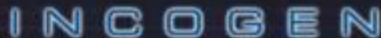
**Slide 23 - Slide 23**



**Slide notes**

An example of a heuristic algorithm is the MSA algorithm, which uses the Carillo-Lipman Bound procedure. This procedure provides a bound in the form of a polyhedron around the diagonal in a hypercube. This bounds the search space for finding the MSA of a set of sequences.

**Slide notes**

First, we consider the pairwise alignments of each pair of sequences and create a phylogenetic tree from these scores.  We then produce a "draft" MSA built from the phylogenetic tree.  The pairwise alignments and draft MSA provide us with a reduced solution space on which we can use dynamic programming to find the solution.  This method does not guarantee an optimal alignment for all the sequences in the group because so much of the solution space is excluded from the search.  It does, however, give us an optimal alignment from within the search space.

**Slide 25 - Progressive Methods**



**Slide notes**

Progressive methods are another way of calculating MSAs. The first steps used in progressive methods are very similar to those used in the MSA algorithm. However, the progressive methods do not, by default, refine the draft MSA by doing a full search in the smaller search area. It also does not guarantee an optimal alignment.

**Slide 26 - Progressive Methods Problems**



**Slide notes**

The progressive methods have a few drawbacks.  First, they are highly sensitive to the choice of initial aligned pairs.  These initial pairs are "frozen" even if subsequent steps show that they are not correct.  For example, there are a couple of equivalent looking pairwise alignments for these two sequences, and the algorithm happens to have chosen this one.  This alignment is now frozen and cannot be changed.  However, as we look at other sequences that we need to align, it becomes obvious that the gap in sequence y should have been placed elsewhere.  The choice of scoring matrices and gap penalties can make it more likely that you will get the correct initial alignments, but choosing the best one is not straightforward.  The likelihood of large errors in the initial alignments increases as the sequences become more distantly related.

**Slide 27 - Progressive Methods Iterative Refinement**



**Slide notes**

Some researchers have attempted to improve upon the progressive methods by adding refinement steps to the procedure.  In this case, we generate the initial alignment and then remove a sequence from the alignment and align it with the remainder of the MSA.  We continue to do this until the alignment score no longer increases.  This procedure is guaranteed to converge to a local maximum.  If the initial alignment is poor, it may not converge to the global maximum.

**Slide 28 - Profile Alignment**



**Slide notes**

Once an alignment has been generated, it is helpful to use the position specific information from the MSA when adding new sequences to the alignment. Essentially, we perform a pairwise sequence alignment of the new sequence to the MSA's profile. Many progressive alignment methods, such as ClustalW, use pairwise alignment of sequences to profiles and profiles to profiles.
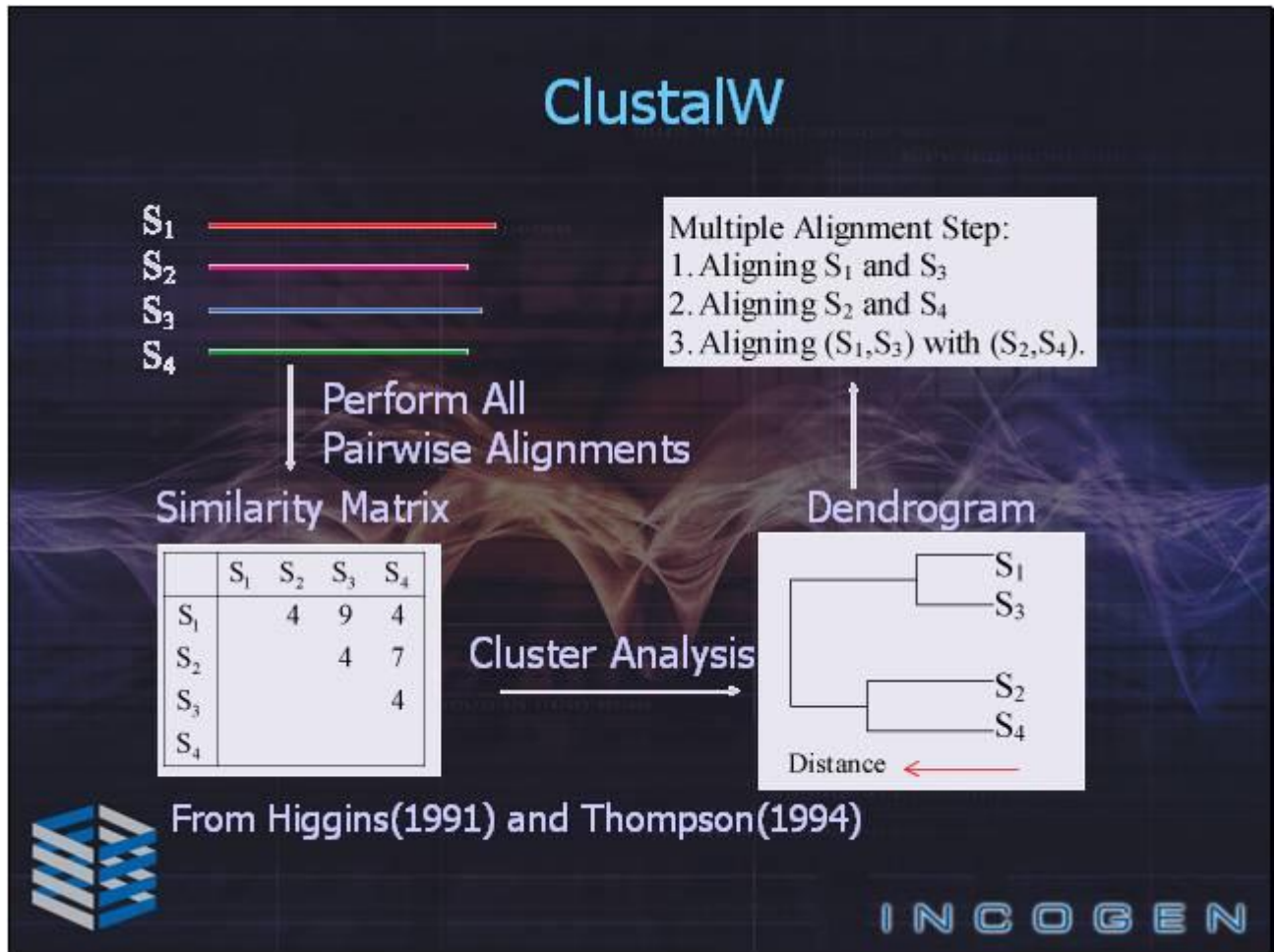
**Slide 29 - ClustalW**



**Slide notes**

ClustalW is currently the most popular multiple sequence alignment algorithm. It performs a pairwise sequence alignment between all of the sequences and constructs a phylogenetic tree from the results. It then combines the results starting from the most closely related groups to the most distantly related groups, using dynamic programming to align the most closely related pairs of sequences.

**Slide notes**

So, ClustalW takes a group of sequences and performs all pairwise alignments.  It then calculates a similarity matrix, which it analyzes to see how distantly related the groups of sequences are.  It then aligns the sequences and groups of sequences, aligning the most closely related at each step, until the MSA is complete.

**Slide 31 - Summary**



**Slide notes**

In summary, choosing the best scoring scheme is critical to the creation of a meaningful MSA. Dynamic programming methods, while guaranteeing an optimal alignment, are too computationally expensive to use for even moderate numbers of sequences, although there are heuristics that can be used to reduce the number of calculations. Progressive methods are much less computationally expensive, but they are very sensitive to initial alignments and may not produce good alignments, especially for distantly related sequences. Using an iterative approach improves the alignments produced by progressive methods, but it is still sensitive to the initial alignment. The profile methods allow integration of position-specific information and profile-profile alignments. In general, most computational methods use a large number of heuristics to obtain an optimal alignment.