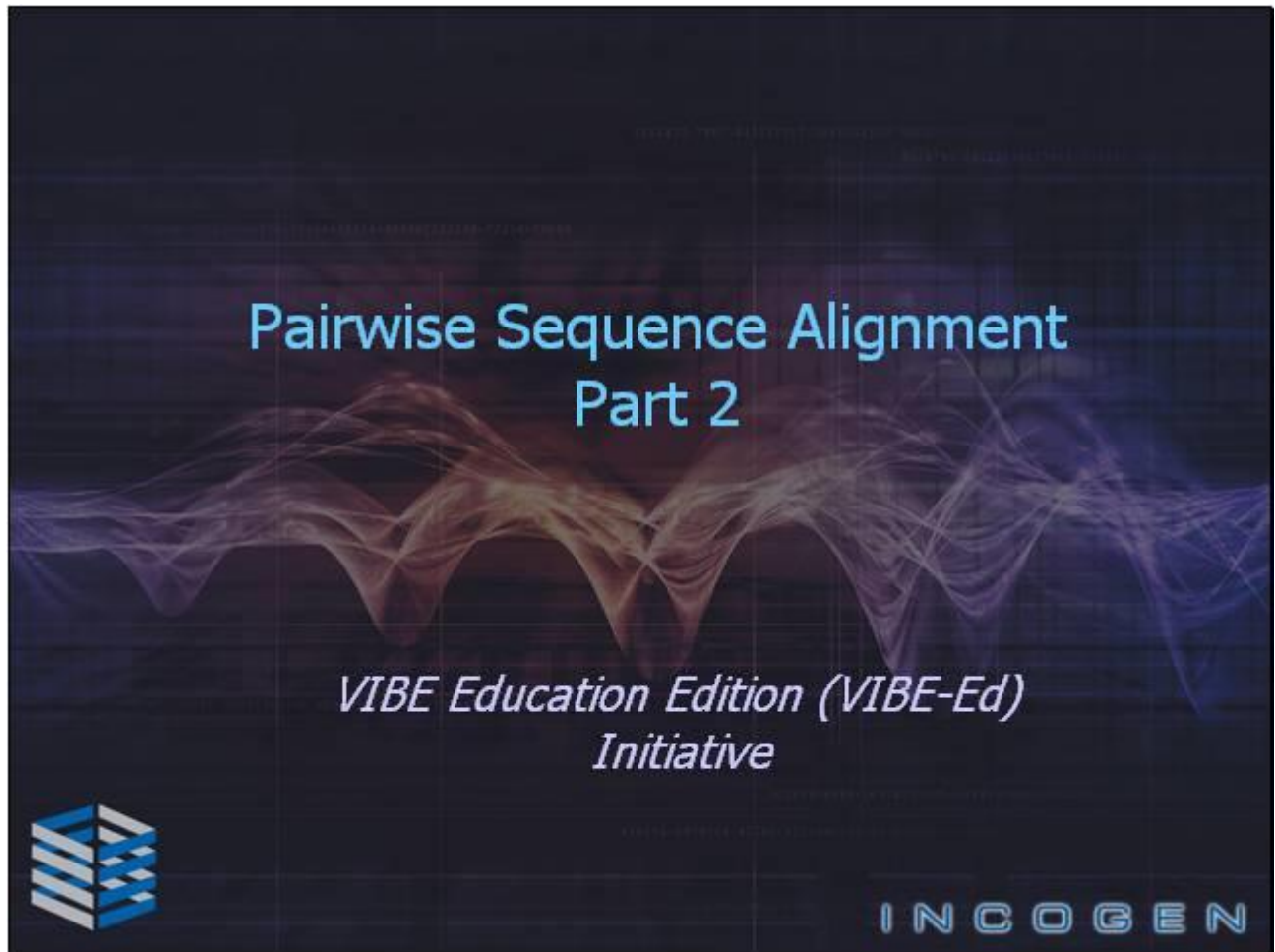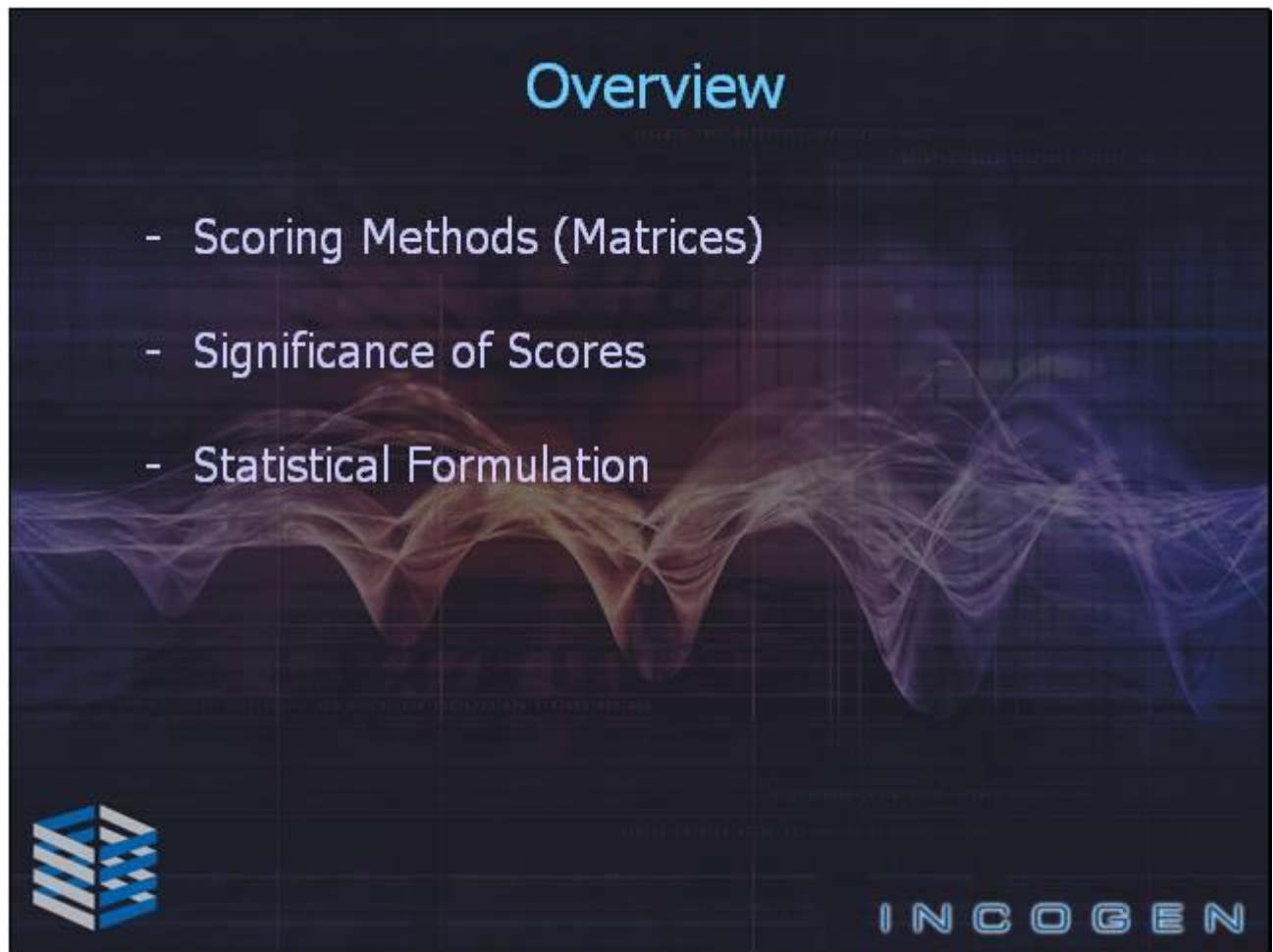**Slide 1 - Pairwise Sequence Alignment Part 2**



**Slide notes**

This presentation continues the discussion of pairwise sequence alignment.

**Slide 2 - Overview**



**Slide notes**

In particular, we are going to have a more in-depth discussion about scoring methods, about the significance of the scores that the algorithms produce, and about the statistics behind the calculations of significance and score.

**Slide 3 - Statistics of Similarity Searches**



## Statistics of Similarity Searches

- Score
    - Using a scoring method, we can generate a maximum scoring alignment.
    - What kind of scoring method should we use?

- Significance
    - How significant is the score?
    - How likely is it that the two sequences that produced the score are related?

INCOGEN

**Slide notes**

Specifically, we want to think about the following questions. First, what kind of scoring method should we use to generate the maximum scoring alignment? Second, how significant is this score; that is, how likely is it that the two sequences producing this alignment score are related?

**Slide 4 - Scoring Methods**



**Slide notes**

First, some more detailed information about scoring methods.

**Slide notes**

To calculate the final score of a potential alignment, each symbol pairing is assigned a numerical value, based on how frequently you would expect to see this pairing in two related species. To simplify calculating the final score, the values must be additive. In general, the values are arranged in a "scoring matrix."

**Slide notes**

The simplest scoring matrix for DNA sequences assigns a value of 1 for a match and 0 for a mismatch.  Using this matrix, the example alignment shown has a score of 5.

**Slide 7 - Slide 7**



**Slide notes**

Using this scoring matrix makes it more difficult to tell the difference between a very good alignment that is short and a longer alignment that is very poor. So, we can introduce negative scoring values to penalize mismatches. For this scoring matrix, a match gets a score of 5, and a mismatch gets a score of -4. Using this modified matrix, the example alignment has a score of -51.

**Slide 8 - Slide 8**



**Slide notes**

Protein scoring matrices are larger and more complicated looking since there is a different probability for each amino acid pairing. To find the value associated with any amino acid substitution pair, find the two amino acids on the edges of the table and read off the value where the row and column intersect. For this example scoring matrix, let's look at Threonine, whose single-letter code is T, and Glycine, whose single-letter code is G. The value for a T:T match is 5, and the value for a T:G mismatch is -2.

**Slide 9 - Slide 9**



**Slide notes**

Different amino acid pairings have different values because amino acids have different biochemical and physical properties that influence their relative substitution rate in evolution. In this image, the amino acids are grouped according to the chemistry of their side groups.
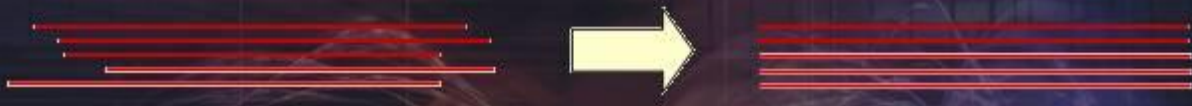
**Slide 10 - Dayhoff PAM Matrix (Point Accepted Mutation)**



### Slide notes

There has been a lot of work to develop meaningful score matrices for protein sequence alignment. The first family of scoring matrices we will look at is the Dayhoff PAM matrix. The matrices list the likelihood that one amino acid could change to another in homologous protein sequences during evolution. It assumes that each amino acid change is independent of any previous changes that occurred at that site. This matrix family was derived from global alignments of protein families that are at least 85% identical.

**Slide 11 - PAM Matrix cont'd**



**Slide notes**

This first matrix, which represents 85% identity, is called PAM 1. PAM 1 reflects an average change of 1% of all amino acid positions. Since we are generally looking at alignments between sequences that are less than 85% similar, other PAM matrices, such as PAM 250, were calculated by multiplying PAM 1 by itself for some number of times; the number of times is reflected in the name of the matrix. Larger PAM numbers mean that that the matrix represents a larger evolutionary distance.

**Slide 12 - Slide 12**



**Slide notes**

Here is the PAM 250 matrix.  This matrix is symmetric, meaning that the value in the matrix is independent of which amino acid out of a pair substituted for the other.

Let's look at Tryptophan, represented by its single-letter code, W, and Cysteine, represented by its single-letter code, C.  The value associated with a W:W pairing, 17, is a very large positive number, while almost all of the other pairings with W, such as W:C, have relatively large negative scores.  Therefore, the matrix represents the fact that Tryptophan is a conserved residue.

**Slide 13 - BLOSUM Matrix  (Blocks Amino Acid Substitution)**



**Slide notes**

The second family of scoring matrices is the BLOSUM matrix.  The BLOSUM matrices were based on the amino acid substitution rate in highly conserved blocks.  The BLOSUM matrices are not explicitly based on an evolutionary model, as the PAM matrices are, but are based instead on related families of proteins.

**Slide 14 - Slide 14**



**Slide notes**

Specifically, the BLOSUM matrices were derived from alignments of domains in distantly related proteins. Occurrences of each amino acid pair in each column of each block were counted, and the totals were used to compute the BLOSUM matrices.

**Slide 15 - Slide 15**



**Slide notes**

The sequences within blocks were clustered based on their level of identity, and the different BLOSUM matrices differ in the percentage of sequence identity used in the clustering.  This percentage is reflected in the matrix name.  Unlike the PAM matrices, larger numbers mean a smaller evolutionary distance.

**Slide notes**

So, which matrix do we choose for our alignment? Generally, BLOSUM matrices perform better than PAM matrices for local similarity searches. It is very important that you remember that large PAM numbers or small BLOSUM numbers should be used for distantly related proteins, while small PAM numbers or large BLOSUM numbers should be used for closely related proteins. The most commonly used matrix for calculating local alignments is BLOSUM 62.

**Slide 17 - Gapped Alignment**



**Slide notes**

Gaps are scored independently of the scoring matrix used.  There are two methods of scoring gaps.  The first simply multiplies the length of the gap by a penalty per length.  The second method uses a gap opening penalty added to a gap extension penalty that is multiplied by the amount the gap is extended (that is, the length of the gap minus 1).

**Slide 18 - Significance of Scores**



**Slide notes**

Now that we know how scores are calculated, let's take a brief look at their significance.

**Slide 19 - Significance of Scores, con't**



**Slide notes**

It is always possible to find an optimal alignment, or Maximum Segment Pair (MSP), between any two protein sequences, regardless of whether they are related. However, not all MSPs are significant. To determine if the MSP is significant, we need to find out how many MSPs with at least that same score we can expect by chance.

**Slide 20 - Two Assumptions**



**Slide notes**

To calculate this, we will make two assumptions.  First, at least one of the target frequencies is positive.  Second, the expected score for aligning a pair of random sequences is negative.

**Slide 21 - Statistical Significance:  Expectation Value**



**Slide notes**

With these assumptions, we can use the extreme value distribution (EVD) to calculate the number of alignments between random sequences that we expect to see with a given score or better.  We define this quantity to be the e-value.  In general, as the score increases, the e-value decreases, indicating that this alignment is more likely to be biologically significant.

**Slide 22 - Statistical Significance: P-Value (probability)**



**Slide notes**

We can also calculate the probability of finding at least one alignment with a score greater than or equal to some given score. We define this quantity to be the p-value, and it can be also used to help determine if the alignment that has been returned is biologically significant. Like the e-value, as the p-value decreases, it is more likely that the alignment is biologically significant.

**Slide 23 - Score, E-value and P-value compared**



## Score, E-value and P-value compared

| Score | E-value $E(S) = K\,m\,n\,e^{-\lambda S}$ | Probability $P = 1 - e^{-E(S)}$ |
|---|---|---|
| 39 | 12 | 0.999995 |
| 41 | 2.9 | 0.947456 |
| 42 | 1.4 | 0.764656 |
| 46 | 0.0842 | 0.080741 |
| 49 | 0.0100 | 0.009925 |
| 52 | 0.0012 | 0.00118 |
| 55 | 0.0001 | 0.0001 |

m = 980, n = 10,030,834,086   (m*n ~$10^{13}$)
K = 1.37, $\lambda$ = 0.711

INCOGEN

**Slide notes**

This table gives a brief comparison of score, e-value, and p-value.  Holding all other variables constant, as the score increases, both the p-value and the e-value decrease.

**Slide 24 - Statistical Significance**



**Slide notes**

You may be asking why not use raw score instead of e-value or p-value since the three are related and the raw score doesn't need any additional calculations.  Well, you can only directly compare the raw scores of alignments if the alignments were created using the same scoring method and against the same search space, which is the product of the database and query sequence.  Therefore, the raw score is not useful in general.  E-values and P-values take the search space and scoring method into account, which means that they can be directly compared to find the best alignment.

**Slide 25 - Normalized (bit) score**



**Slide notes**

Another useful value used to measure the strength of a pairwise alignment is the bit score, S', derived in this slide.

**Slide 26 - Statistical Formulation**



**Slide notes**

And now, a few words about the statistics behind these models.

**Slide 27 - Statistical Formulation**



**Slide notes**

In order for the scores to be useful, we need to make sure that the scoring methods we use are statistically valid.

**Slide 28 - Probability/Statistics 101**



**Slide notes**

Let's go over a few terms used frequently in statistics. A model is a system that simulates the object under consideration. A probabilistic model is a model that produces different outcomes based on a set of probabilities. In our case, the objects we are simulating are sequences, and the model is a family of related sequences.

## Slide notes

Let's look at a more general example-that of a six-sided die. There are 6 outcomes for each roll of a die, each with their own probability. For a "normal" die, the probability of any given number appearing after a roll is the same: 1/6. If the die is "loaded," the probabilities of the six outcomes is not the same. Regardless of whether the die is "normal" or "loaded," the probability of any outcome must be greater than or equal to zero and the sums of the probabilities of all of the outcomes must equal 1. If the events we are modeling are independent of each other, such as rolls of a die, then we can calculate the probability of a sequence of events by multiplying together the probabilities of the events occurring separately.

**Slide 30 - Biological Example**



**Slide notes**

Now, let's look at a biological example of an independent system: a DNA or protein sequence. A sequence is built from an "alphabet" of residues. A DNA sequence's alphabet contains 4 letters, while a protein sequence's alphabet contains 20. To create a random sequence, we can assume that each residue occurs with a certain probability, regardless of the other residues in the sequence. This allows us to calculate the probability of a particular sequence by multiplying together the probabilities of all of the residues in the sequence. This is called the "Random Sequence Model."

**Slide 31 - Conditional and Joint Probabilities**



**Slide notes**

Let's go back to our previous example involving a die, except now, we are going to look at a system that contains two dice. Each of these dice could have a different probability of rolling a given number, so we must specify both the die and the outcome of the roll in order to find the probability; this is called the conditional probability. (For example, what is the probability of rolling a 2 given the condition that you are rolling the red die?) If the condition isn't guaranteed, then the probability depends on meeting that condition and then getting the outcome; this is called a joint probability. (For example, what is the probability of rolling a 2 on a red die if you pick from a bag that contains both red and blue dice?) To calculate a joint probability, you multiply the probability of meeting the condition with the probability of getting the outcome once the condition is met.

**Slide 32 - "Occasionally dishonest" casino**



**Slide notes**

Let's look at another concrete example. Let's say, for instance,. That you are in an "occasionally dishonest" casino where 99% of the dice are fair and 1% are "loaded" so that the probability of rolling a 6 is 50%. So, we know that the probability of rolling a 6 on a "loaded" die is 50%, and the probability of rolling a 6 on a "normal" die is approximately 17%. To find the general probability of rolling a 6, we must take into account both the normal dice and the loaded dice. To do this, we multiply the probability of getting a "loaded" die by the probability of rolling a 6 with a "loaded" die (50%). We do the same for the "normal" die, and then sum the two probabilities together.

**Slide 33 - Substitution Matrices**



**Slide notes**

Let's look a little at how this relates to substitution matrices. The general issue: given a pair of sequences, we want to assign a score to an alignment that gives us some idea of the relative likelihood that the two sequences are related versus unrelated. To do that, we need to develop models for the related and unrelated cases that calculate the probability of each case and then take a ratio of the two probabilities.

**Slide 34 - Unrelated (Random) Model R**



## Unrelated (Random) Model R

- Assumes that residue *a* occurs independently with frequency $q_a$
- The probability of the two sequences is just the product of the probabilities of each residue:

$$P(x,y \mid R) = (q_{x1}*q_{x2}*...*q_{xm}) * (q_{y1}*q_{y2}*...*q_{yn})$$

$$= \prod q_{xi} \quad * \quad \prod q_{yi}$$

INCOGEN

**Slide notes**

First, let's look at the model for unrelated sequences. In this model, we assume that each residue occurs independently of any other in the sequence with some given probability. The probability that the two sequences exist independently of each other is just the product of probabilities of each residue in each sequence.

**Slide 35 - Match Model M**



## Match Model M

- Assumes that residue *a* occurs independently with frequency $q_a$
- The probability of the two sequences is just the product of the probabilities of each residue:

$$P(x,y \mid R) = (q_{x1} {}^{*} q_{x2} {}^{*} \ldots {}^{*} q_{xm}) {}^{*} (q_{y1} {}^{*} q_{y2} {}^{*} \ldots {}^{*} q_{yn})$$

$$= \prod q_{xi} {}^{*} \prod q_{yj}$$

INCOGEN

**Slide notes**

Now let's look at the model for related sequences.  In this case, each aligned pair of residues occurs with some joint probability.  This joint probability represents the idea that the two residues are derived from some unknown common ancestor.  The probability that these aligned pairs exist is given by the product of all of the joint probabilities.

**Slide 36 - Odds Ratio**



**Slide notes**

To get these probabilities into a form that can be combined by adding, we need to find the log-odds ratio score. First, we find the odds ratio by dividing the probability that the two sequences are related by the probability that the two sequences exist independently of each other. We then take the log of that to find the log-odd ratio score. Since each aligned pair is independent of the other aligned pairs in the sequences, we can also calculate a log-odds ratio for the aligned pairs independently. It is this value that appears in scoring matrices.

**Slide 37 - Substitution Matrices - Revisited**



**Slide notes**
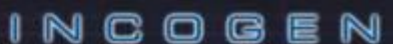
To reiterate, scores are only statistically meaningful when the appropriate scoring matrix is used.

**Slide 38 - Low-complexity Regions**



## Low-complexity Regions

- Significant percentage of regions with highly biased composition
- This is due to:
  - retrotransposons
  - ALU region
  - microsatellites
  - centromeric sequences, telomeric sequences
  - 5' Untranslated Region of ESTs
- Example of EST with simple low complexity regions:

T27311
GGGTGCAGGAATTCGGCACGAGTCTCTCTCTCTCTCTCTCTCTCTCTC
TCTCTCTCTCTCTCTCTCTCTCTCTCTCTCTCTCTCTCTCTCTCTCTC

- Repetitive sequences increase the chance of a high-scoring, but most likely meaningless, alignment during a database search.

INCOGEN

**Slide notes**

The last thing we'll discuss in this lecture is the effect of low-complexity regions on alignments and scoring. Some sequences contain regions that are highly repetitive and low-complexity. This can occur because of a variety of reasons. These repetitive regions increase the chance of a high-scoring, but meaningless, alignment during database searches. Therefore, it is better to mask these regions before beginning a pairwise search.

**Slide 39 - Summary**



**Slide notes**

In summary, you must take speed and sensitivity into account when choosing the appropriate algorithm for determining alignments.  To further speed up the algorithms, use the smallest database that will answer your question. Take some thought when picking which scoring matrix to use; the default matrix may not provide you with meaningful results.  The score of any alignment increases with the size of the search, which is why it is advisable to use p-values and e-values for comparison purposes.  Finally, filter out or mask low-complexity regions in sequences to prevent spurious, but high scoring, alignments from appearing in your results.