# Purpose of Color Schemes

In sequence alignment, it's useful to color
the amino acids in a way that predicts
their similarity, so we can check
with a glance the estimated alignment.

Typically, this is done manually, by
experienced scientists, using colors according
to pre-determined physical properties,
such as charge.

This isn't perfect multiple-fold, as two amino acids
having a similar property doesn't always ensure that
their alignment will be optimal in the given context.
Also, we would like an automated solution for every
situation.

# Substitution Matrix Schema

DNA ANALYSIS

A better way to do this is to use the given substitution matrix and a chosen color space, and define through them a scoring function. Then, optimize that to approximate the optimal.

Given the substitution matrix $M$, define

$$D' \ni \left( D'_{ij} = \begin{cases} \frac{M_{ii} - M_{ij} + M_{jj} - M_{ji}}{2}, & j \leq i \\ 0, & j > i \end{cases} \right) \xrightarrow{\text{scaled to average 1}} D$$

and the chosen color space $CIE\ L*a*b*$, define

$$C \ni \left( C_{ij} \approx \sqrt{(L_i^* - L_j^*)^2 + (a_i^* - a_j^*)^2 + (b_i^* - b_j^*)^2} \right)$$

These are triangular matrices that define the pairwise distance and color variance respectively.

# Substitution Matrix Schema

We can now define the score function

$$S_T = S_H + S_C$$

where

$$S_H = \sum_{ij}(f_s C_{ij} - D_{ij})^2$$
$$S_C = \frac{f_C}{\langle C \rangle}$$
$$f_s = \frac{\langle D \rangle}{\langle C \rangle}$$

$\langle \cdot \rangle$ the arithmetic mean

$f_C$ a given hyperparameter

which essentially defines the total square error of the scaled color variance to distance for each pair, $S_H$, biased by the contrast enabler $S_C$.

This function describes a (generally) non-convex problem that adequately combines the data into the desired effect: coloring based on distance. It can be optimized with Simulated Annealing.
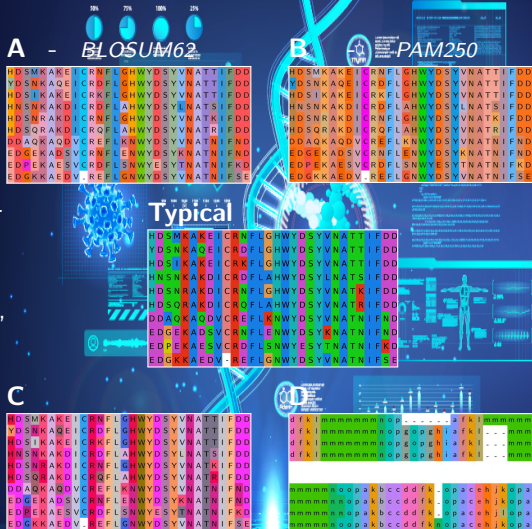
# Results

The figure grid compares this algorithm's outputs (A, B, C and D) to some given color schema made manually.

We can see that the approaches agree on color similarity, but differ on color dissimilarity; the palette on A and B is more gradient and nuanced than the one on 'Typical', which has very sharp changes.

The versatility of the method is also emphasized; Figure C is adapted to red-green colorblindness by removing green as an option. Figure D uses a different alphabet altogether, the protein blocks.



A - *BLOSUM62*

B - *PAM250*

Typical

C

# Conclusions

This tool does a terrific job projecting the relationships between amino acids anchored on molecular evolution rather than anything else, thus tying them to the problem of sequence alignment.

It enables the automatic and independent calculation of sequence alignment heat maps, without the need for research of physical characteristics and unnecessary modeling.

Its generic definition also allows it to be applied to fields other than bioscience, be it for understanding exotic alphabets, cryptographic concepts and much more.

# References

Principal paper:

○ Kunzmann, P.; Mayer, B.E.;
  Hamacher, K. (2020)
  *"Substitution matrix based color schemes
  for sequence alignment visualization"*

Secondary sources:

○ The course's slides and notes.
○ Wikipedia   *"Multiple sequence alignment"*