

Similarity Search with DNA

Using DNA to do things silicon computers have traditionally done

Overview

What is similarity search?

Why would you want to use DNA?

How do you use DNA?

How could we make this better?

How do you use Cas9 to perform similarity search?

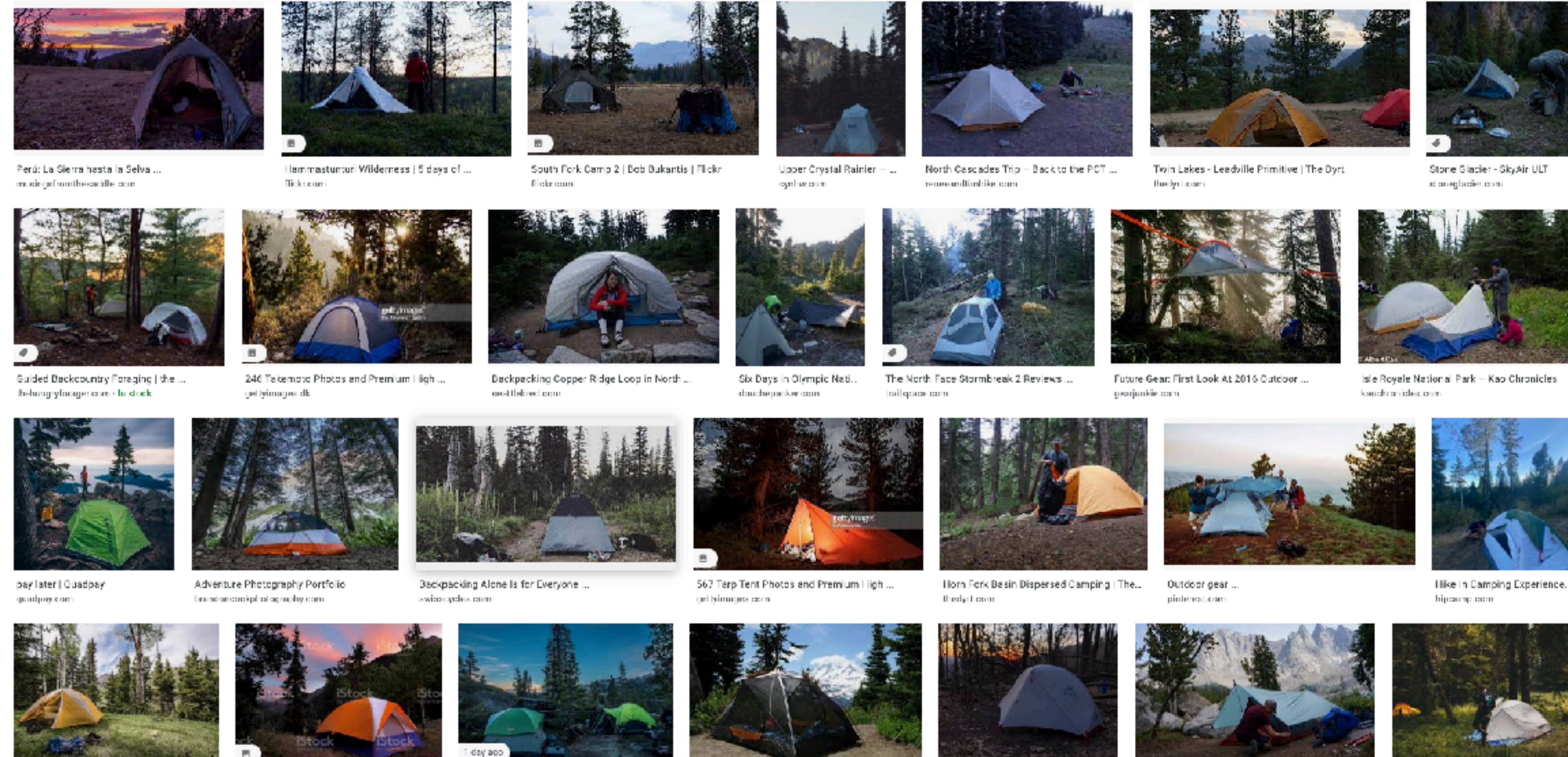
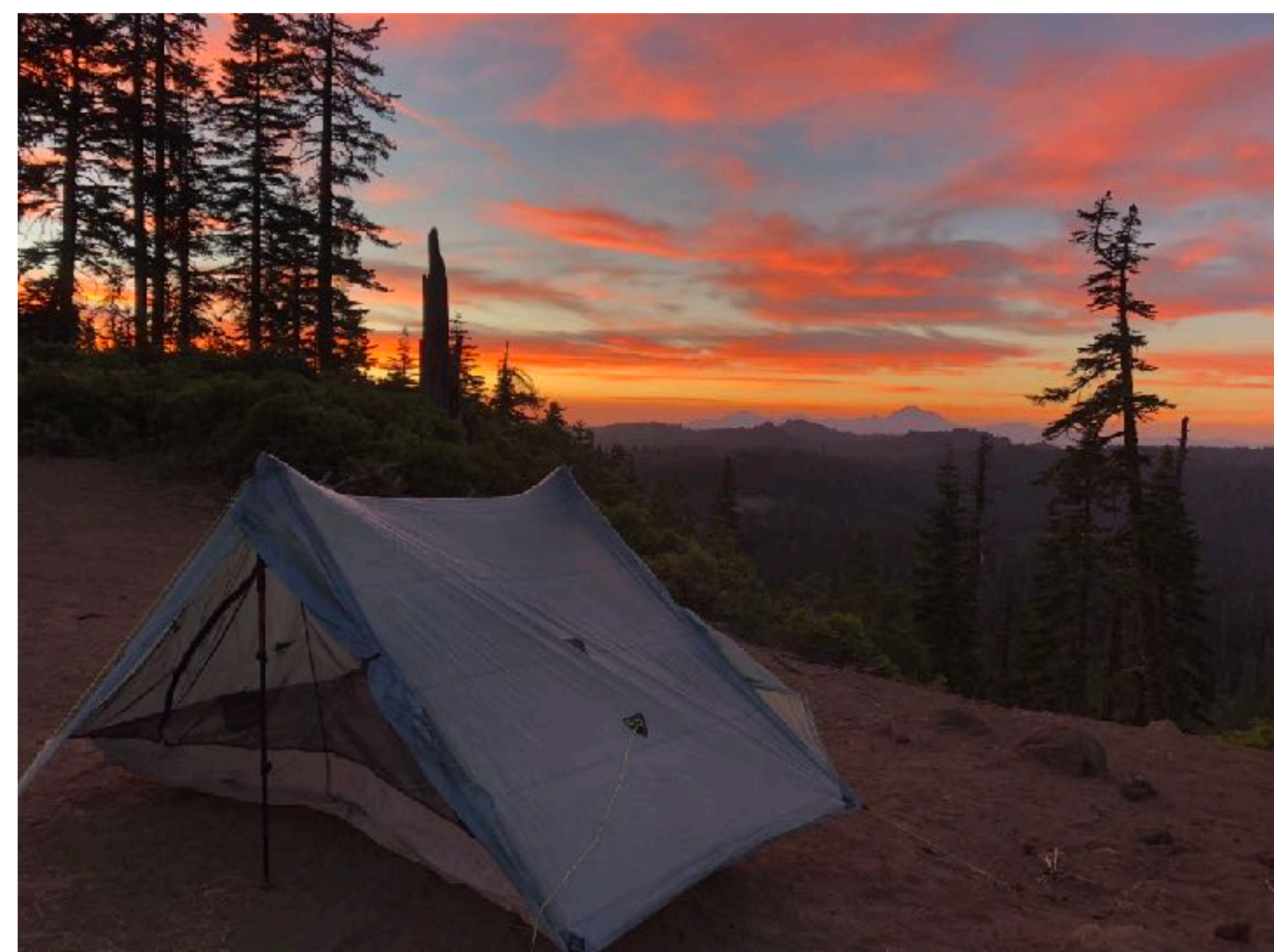
Musings on the intersection of Comp Bio, Molecular Computing, SynBio

What is similarity search?

Input

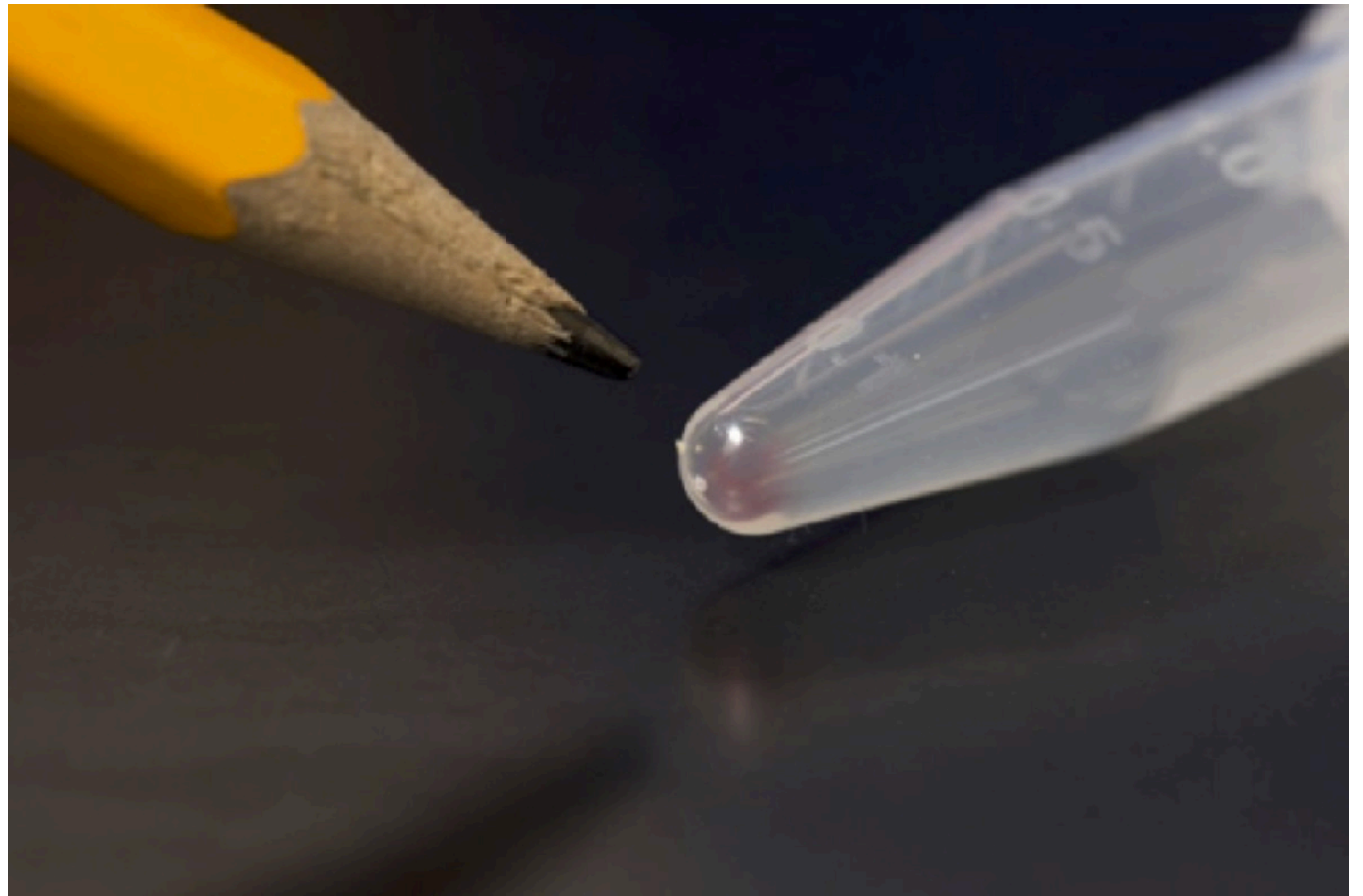
Output

Similar Images



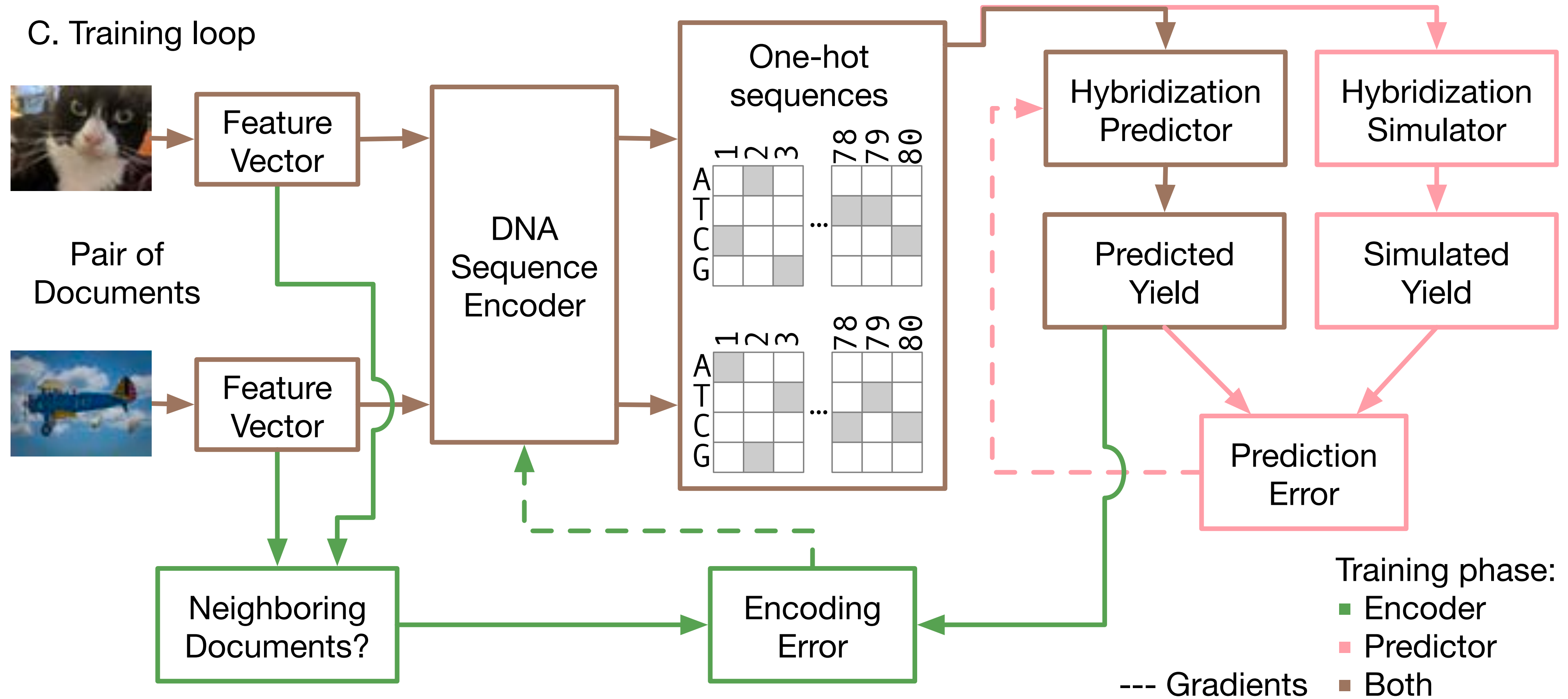
Why would you want to use DNA?

- Parallelism
- DNA information density
- DNA longevity
- Ease of distributing DNA databases
- Sometimes faster
- Sometimes more energy efficient



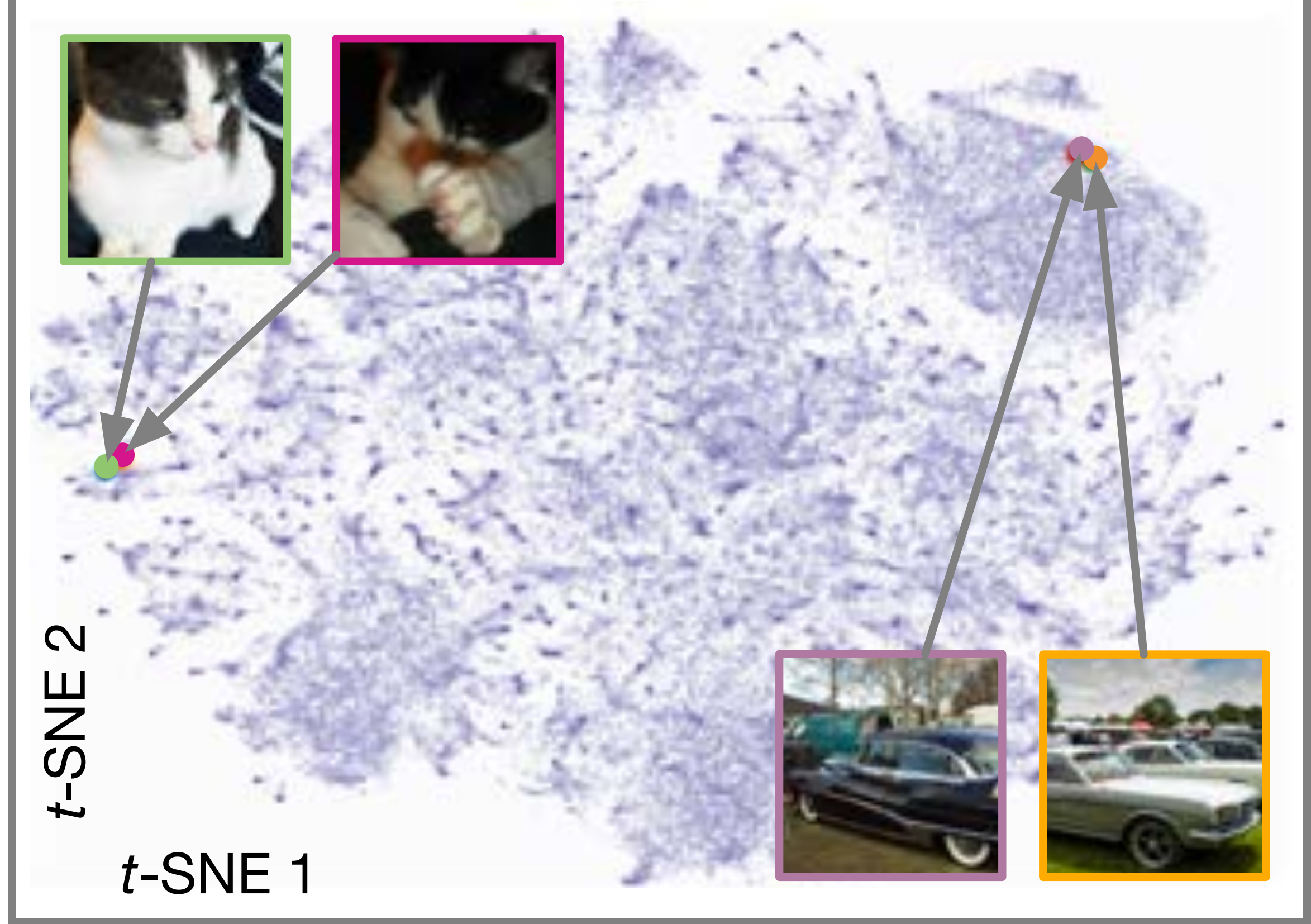
In that faint pink smear is ~10TB of data

Similarity Search



Similarity Search

A. Document similarity as geometric space



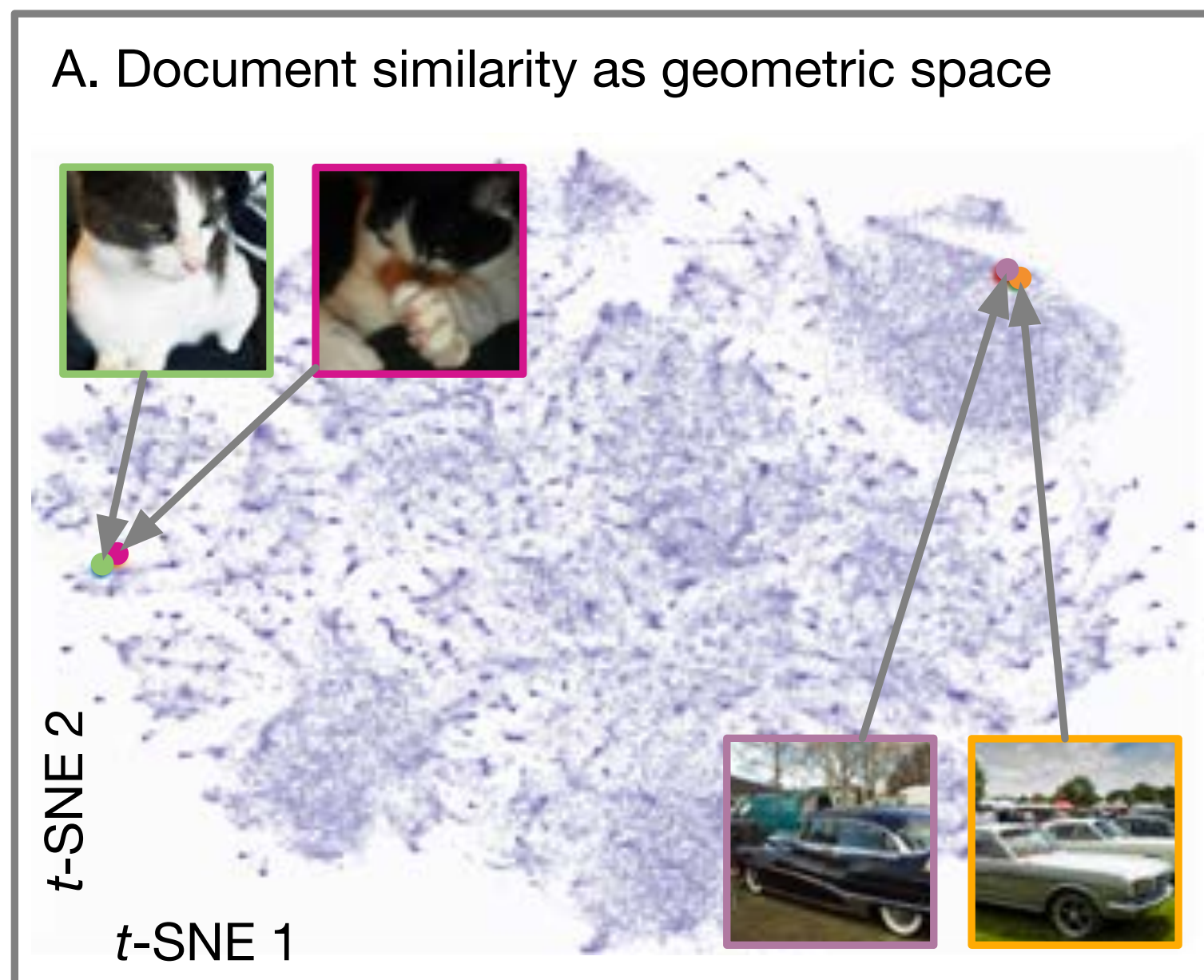
Curse of high dimensionality:

Exact indexing schemes in high-dimensionality spaces are no better than a costly linear or “brute force” search, which is infeasible for large databases

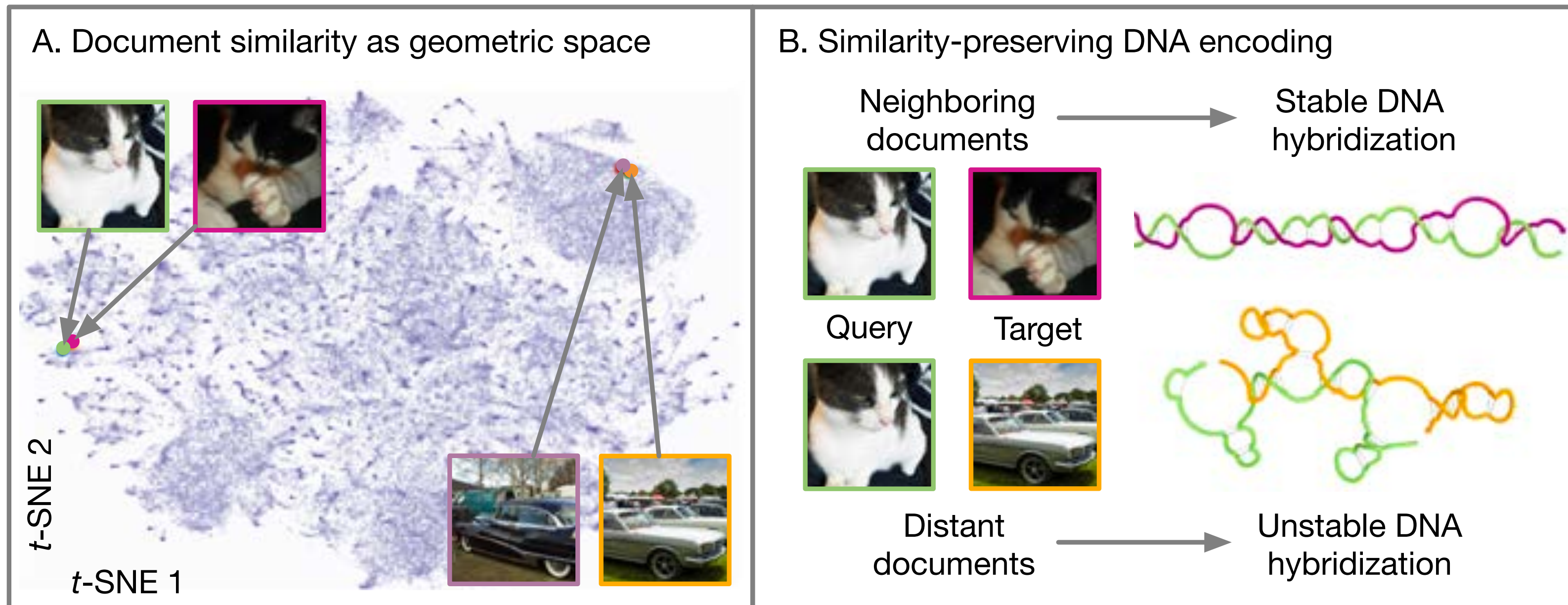
Tradeoff:

Rather than finding exact nearest neighbors, our goal is to maximize the number of near neighbors retrieved while minimizing the number of irrelevant results

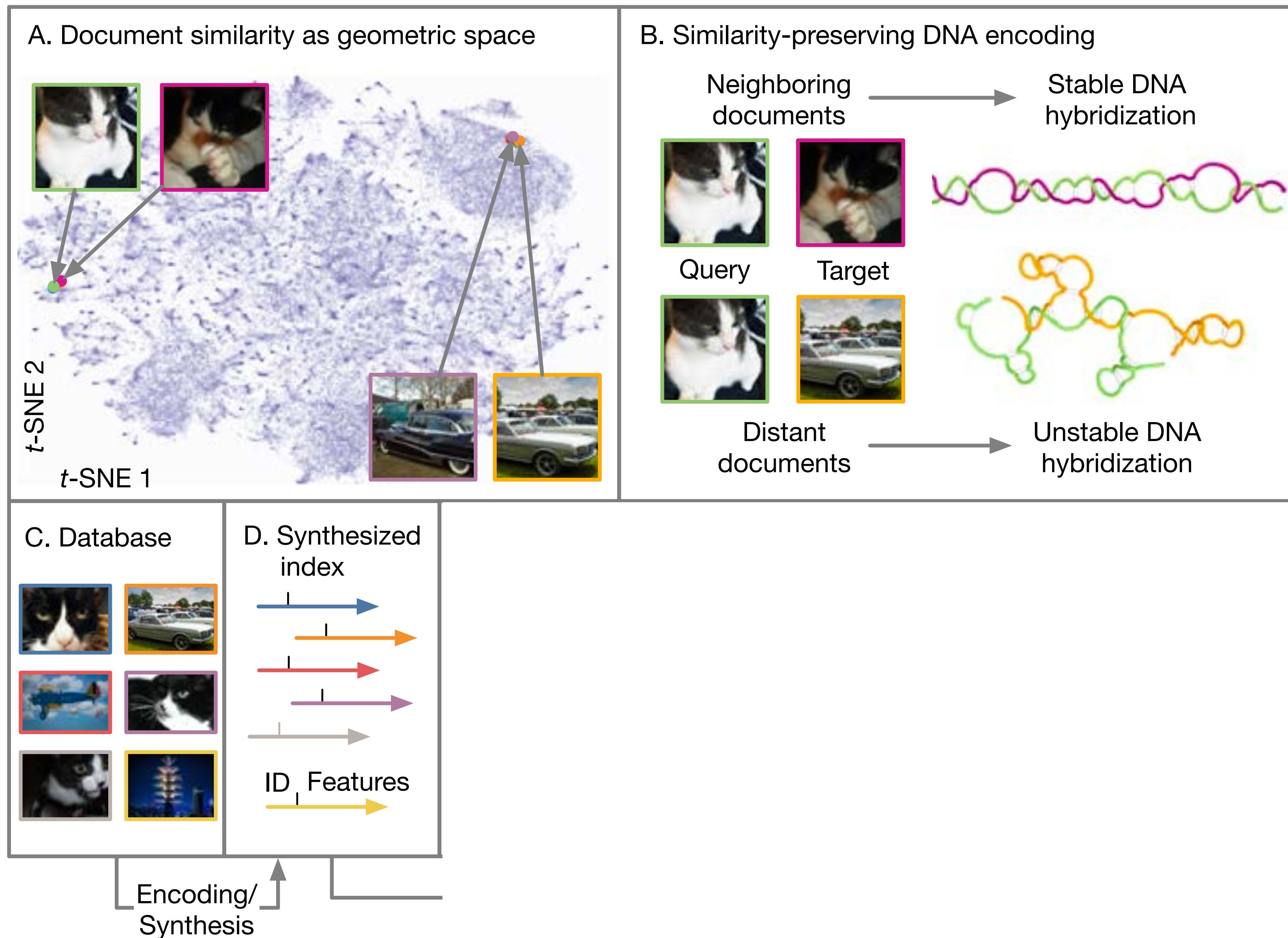
How do we do similarity search with DNA ?



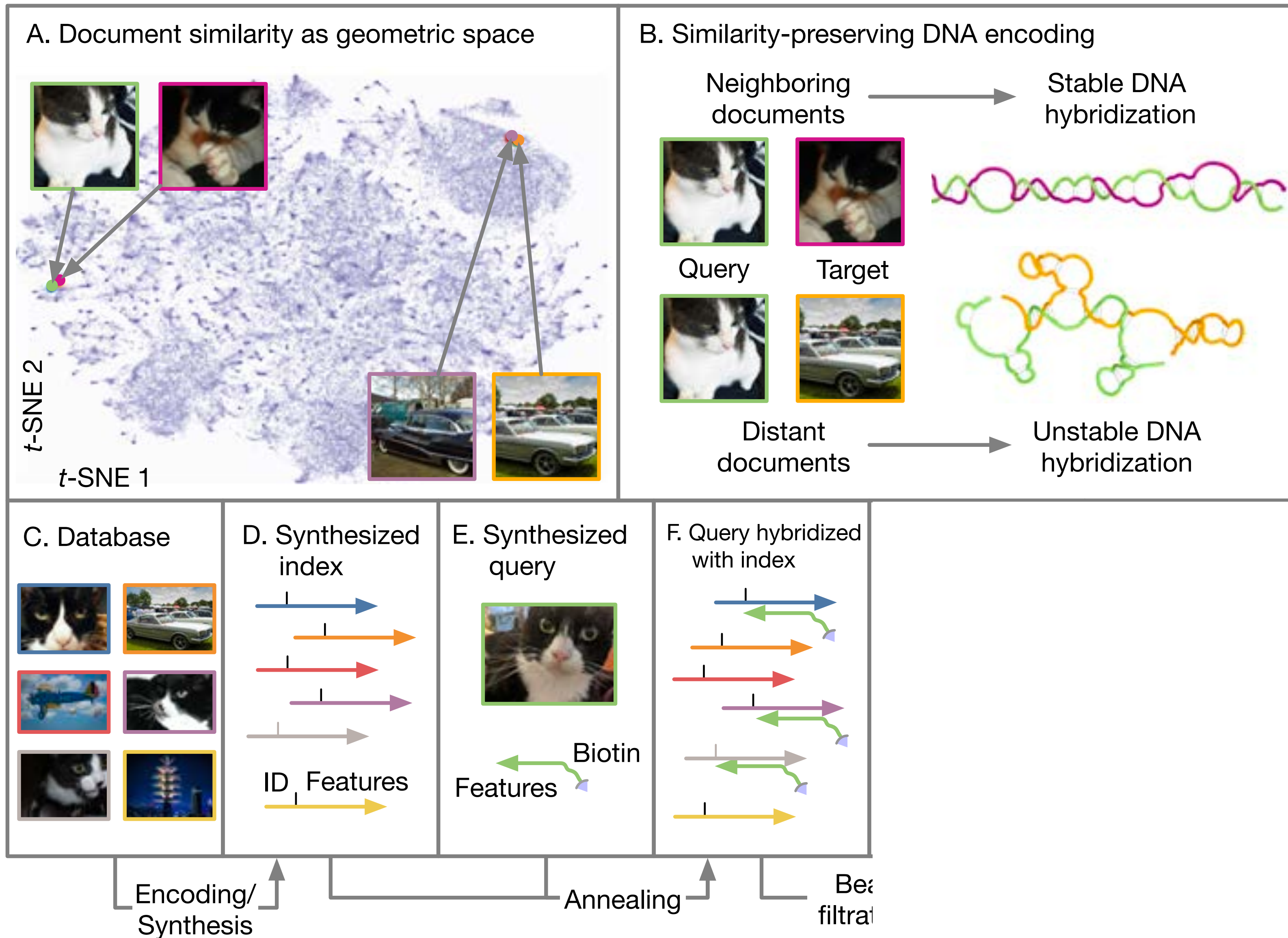
How do we do similarity search with DNA ?



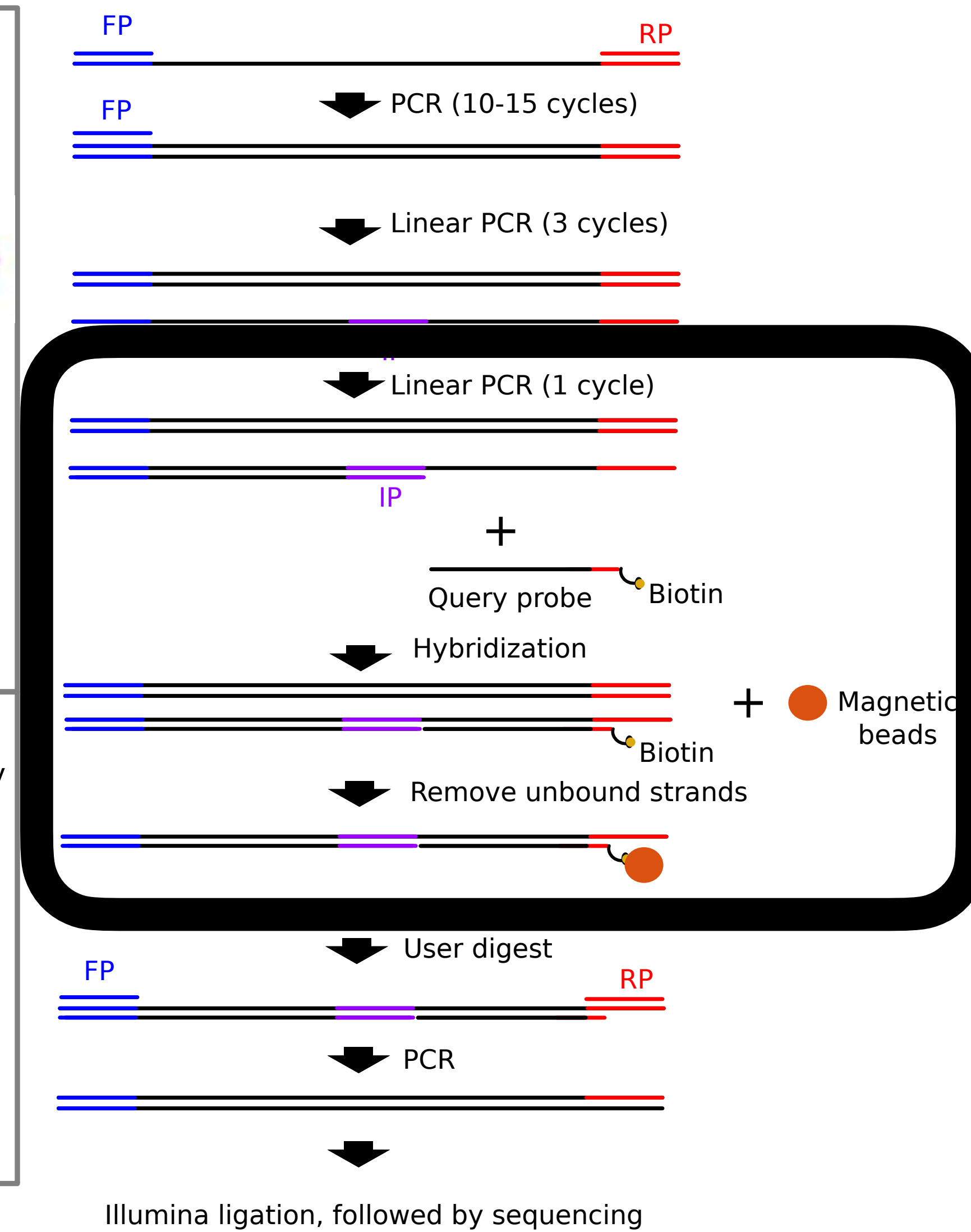
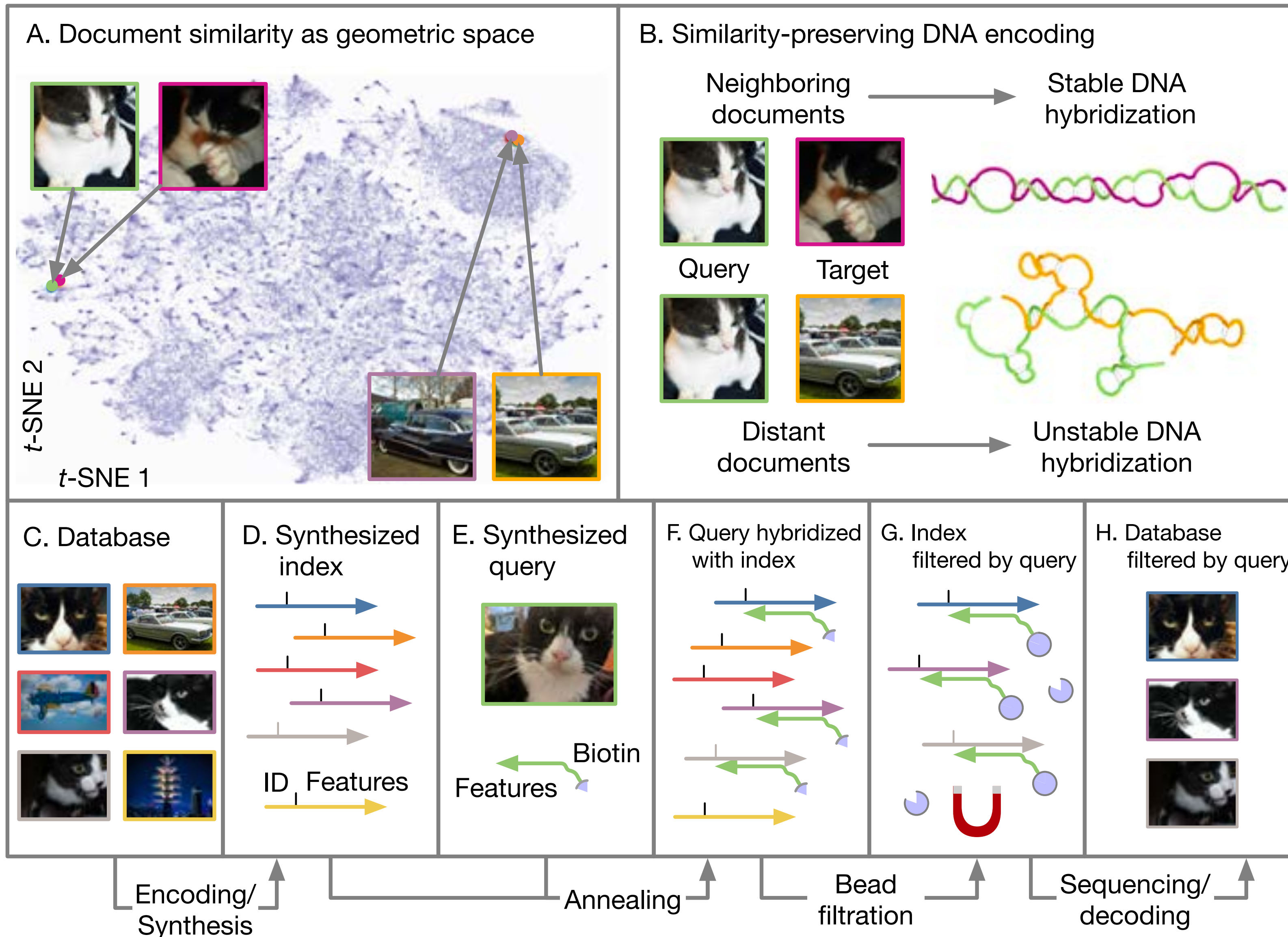
How do we do similarity search with DNA ?



How do we do similarity search with DNA ?



How do we do similarity search with DNA?

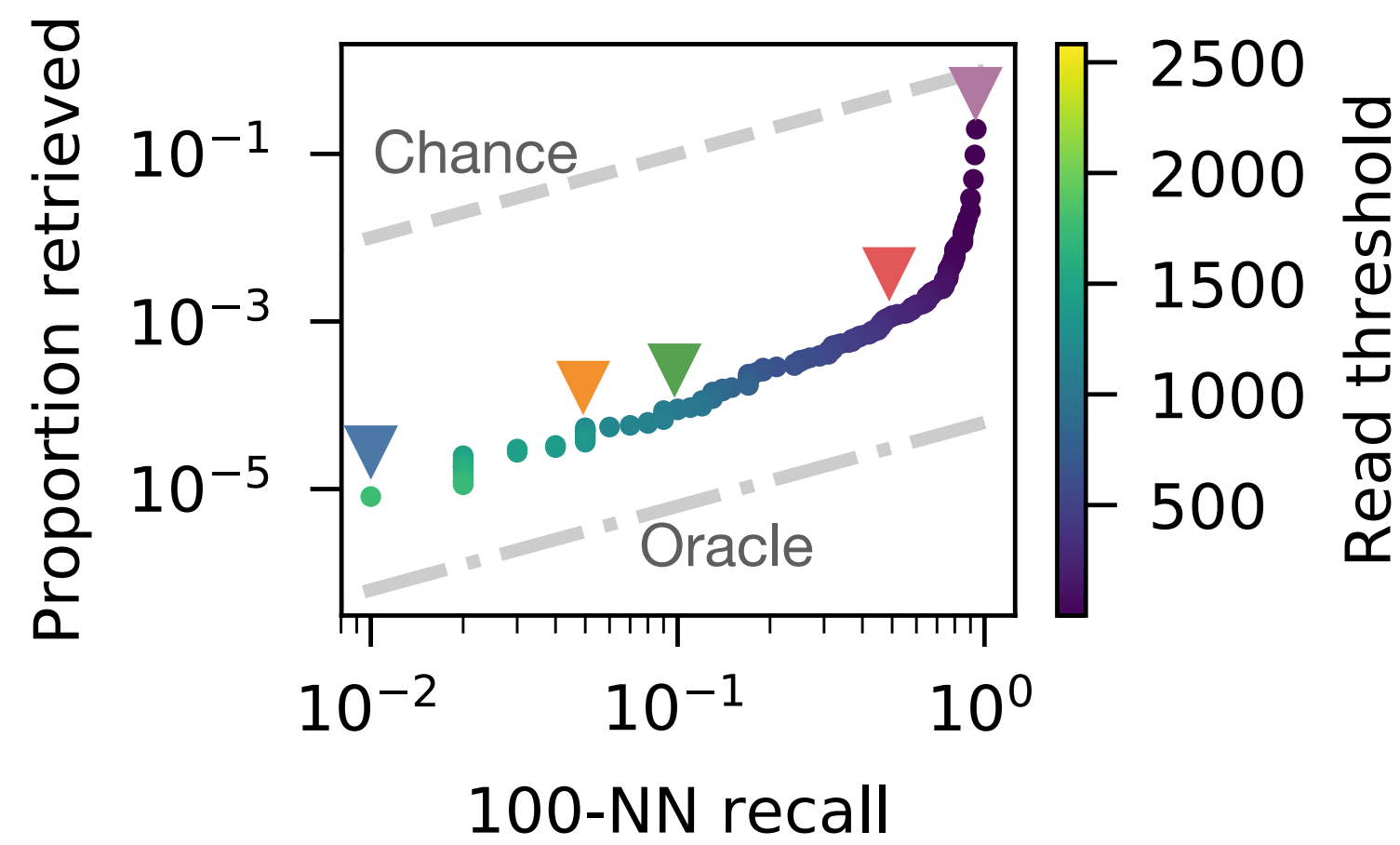
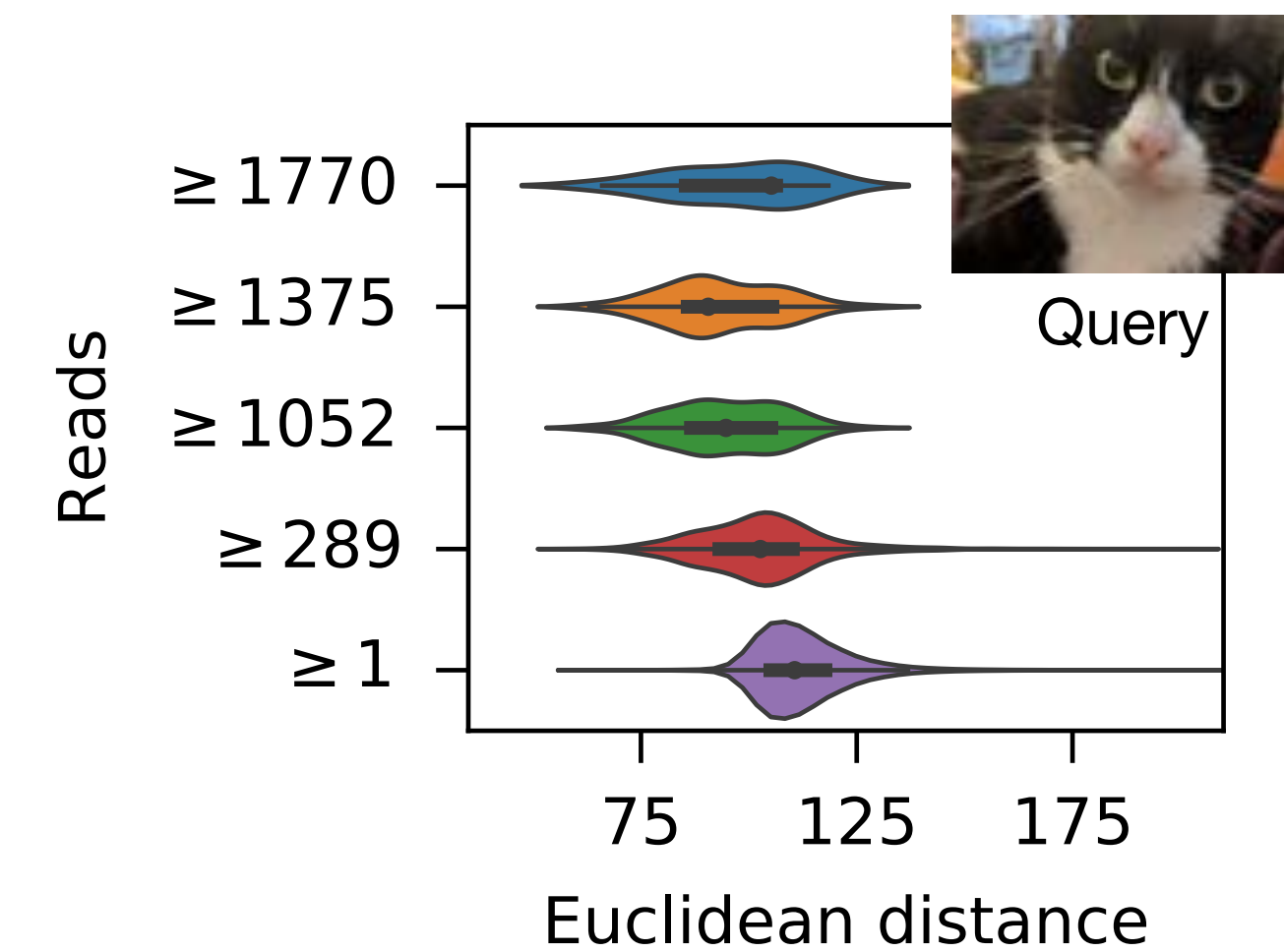


It works!

A. Distribution of similarity across read depths

B. Retrieval as a function of read depth

C. Sets of retrieved images for select read depth thresholds



The proportion of the entire dataset that must be retrieved (y-axis) to retrieve a certain proportion of the 100 most similar images (x-axis)

	Number retrieved	100-NN recall	Top 5 results (sorted by similarity)				
▶	13	0.01					
▶	58	0.05					
▶	141	0.10					
▶	1,831	0.50					
▶	315,736	0.94					

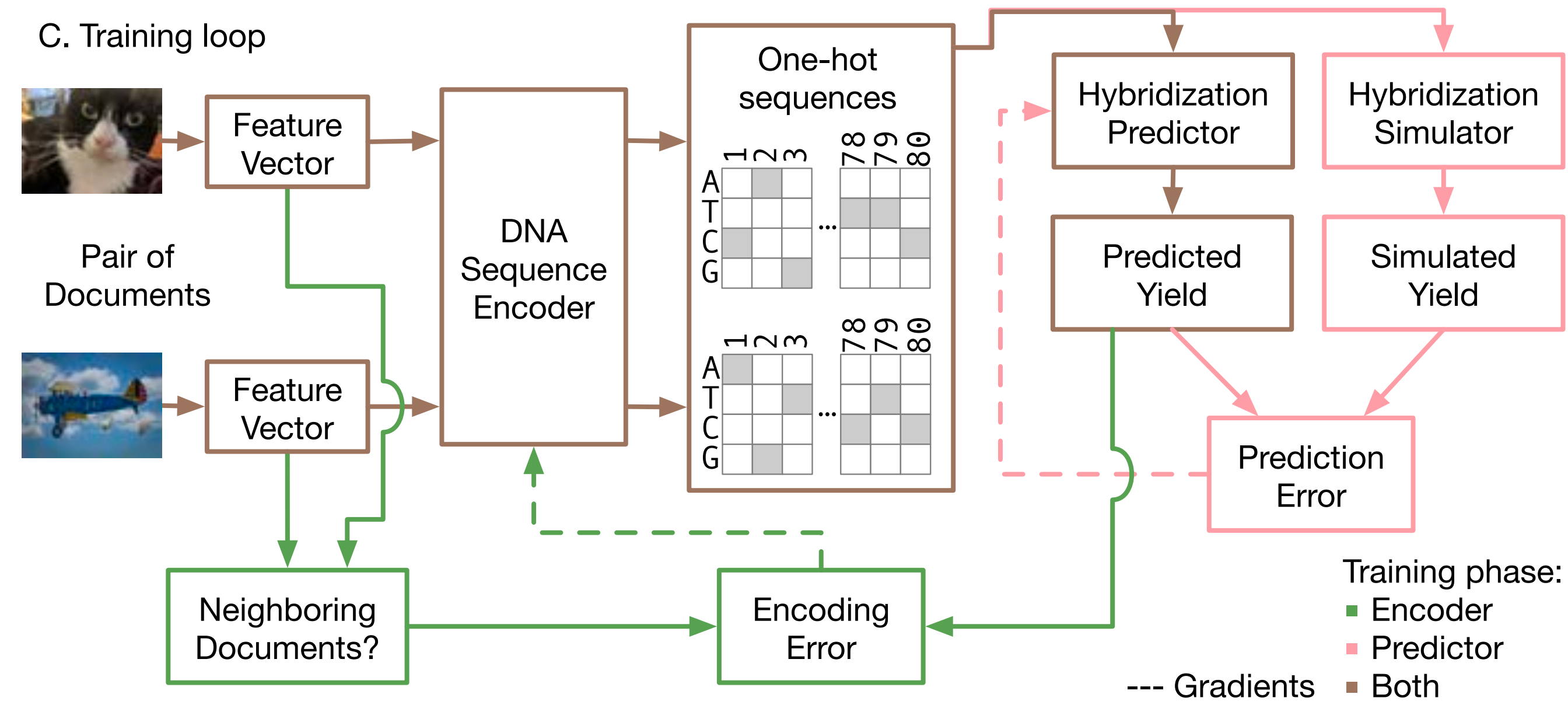
How can we improve?

- More energy efficient
- Faster

...what if we used Cas9?

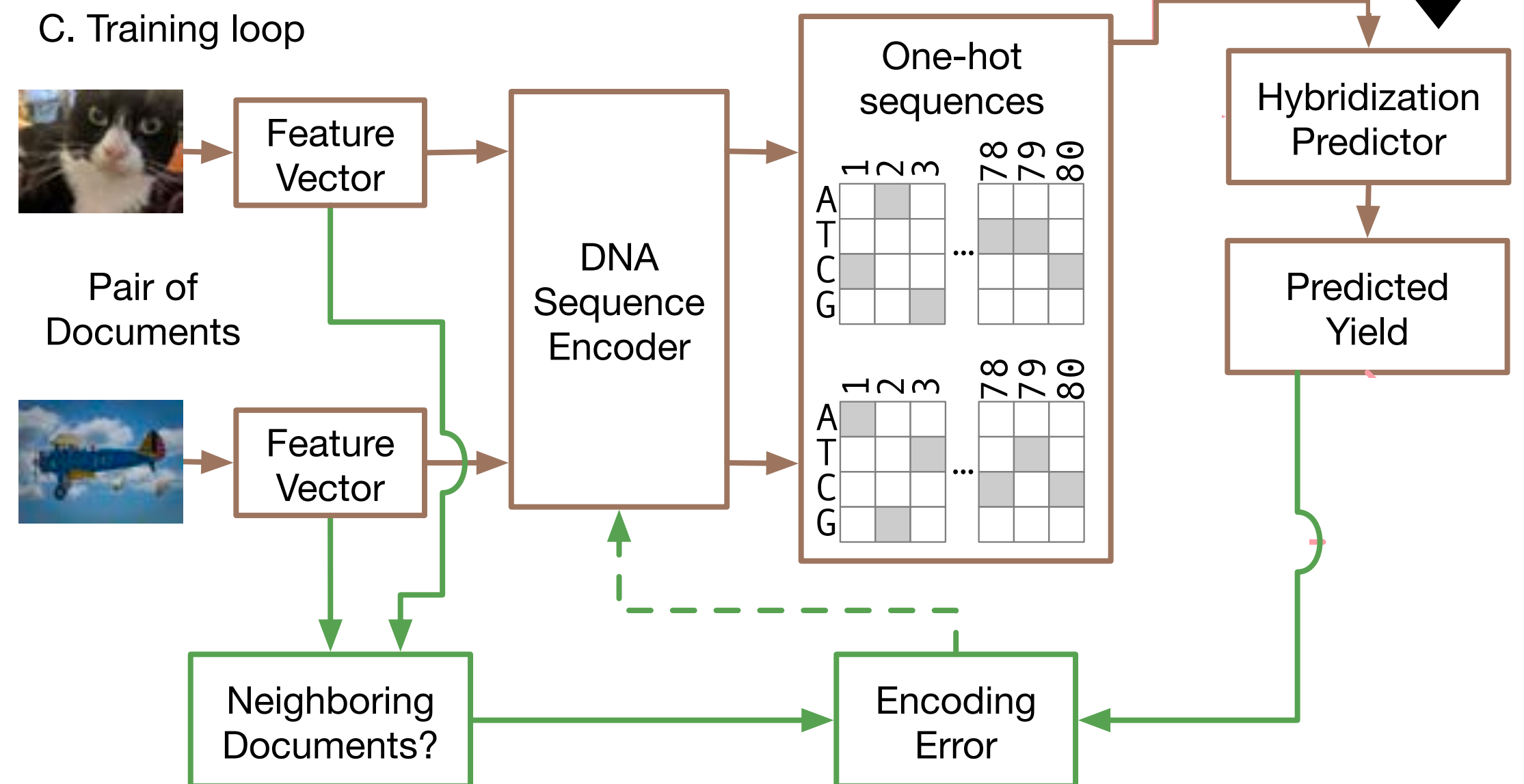
Using Cas9

Hybridization and Bead Extraction



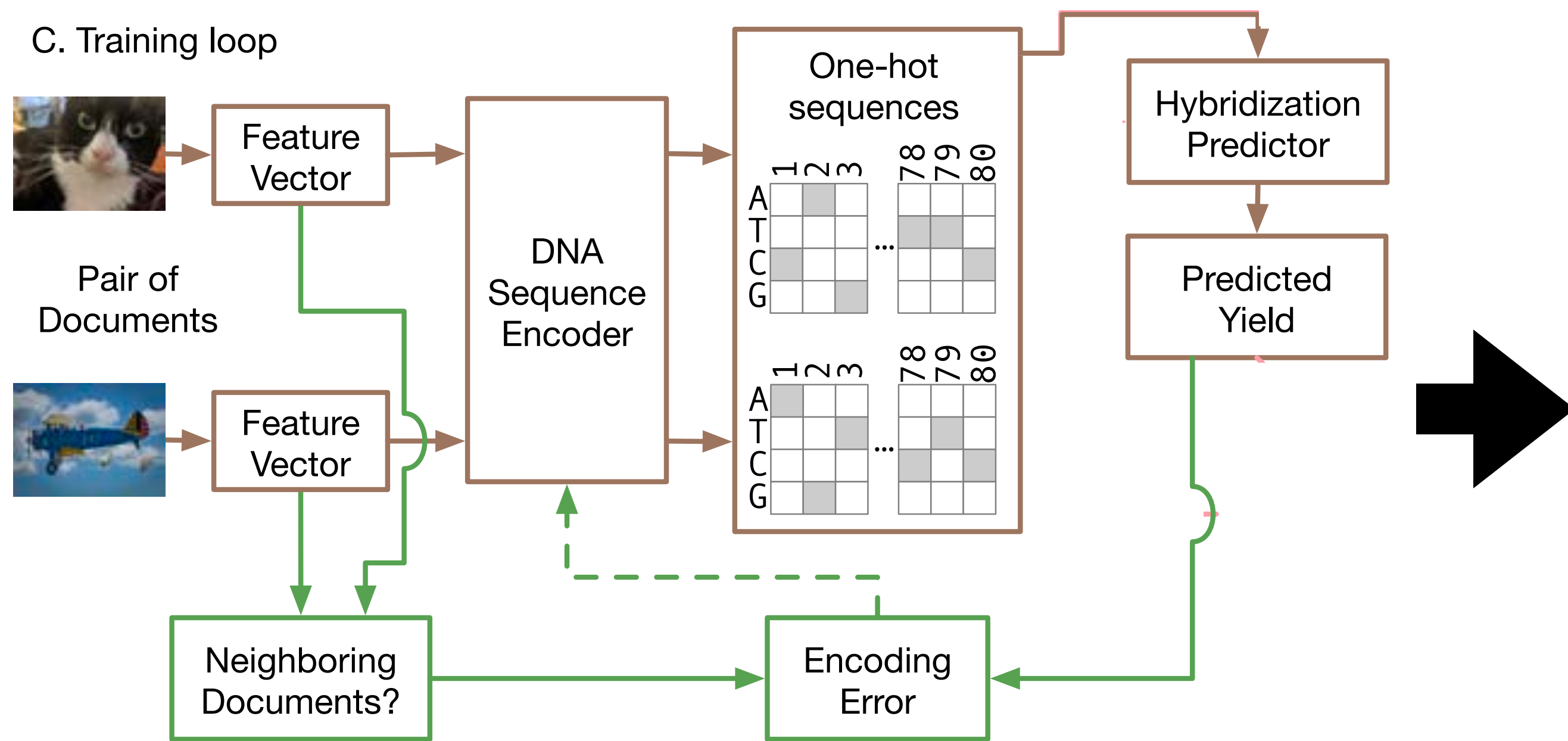
Cas9

*Cas9 Binding Predictor



Using Cas9

C. Training loop



Cas9



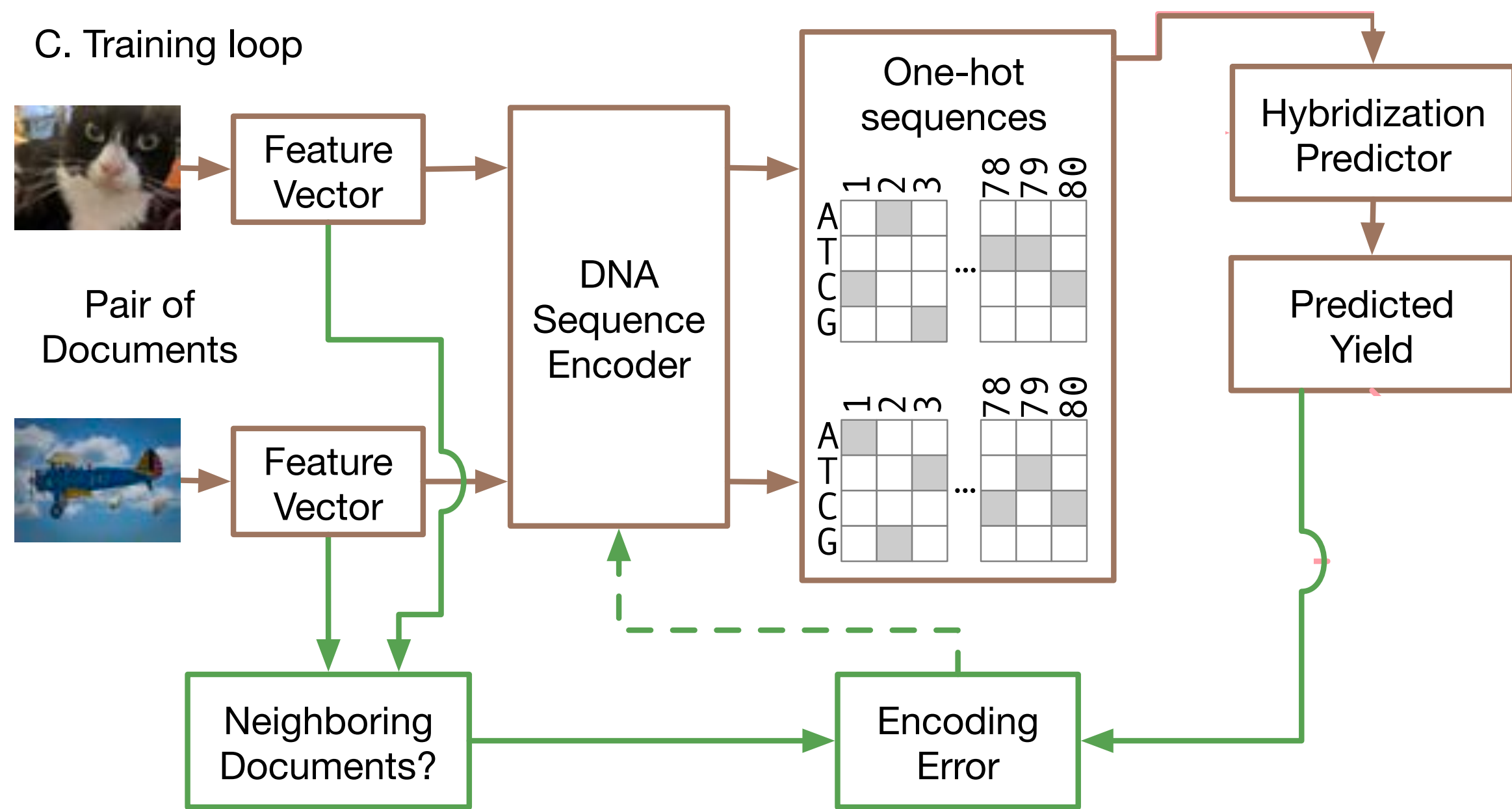
20 features Data ID

We sequence every strand that's cut by Cas9



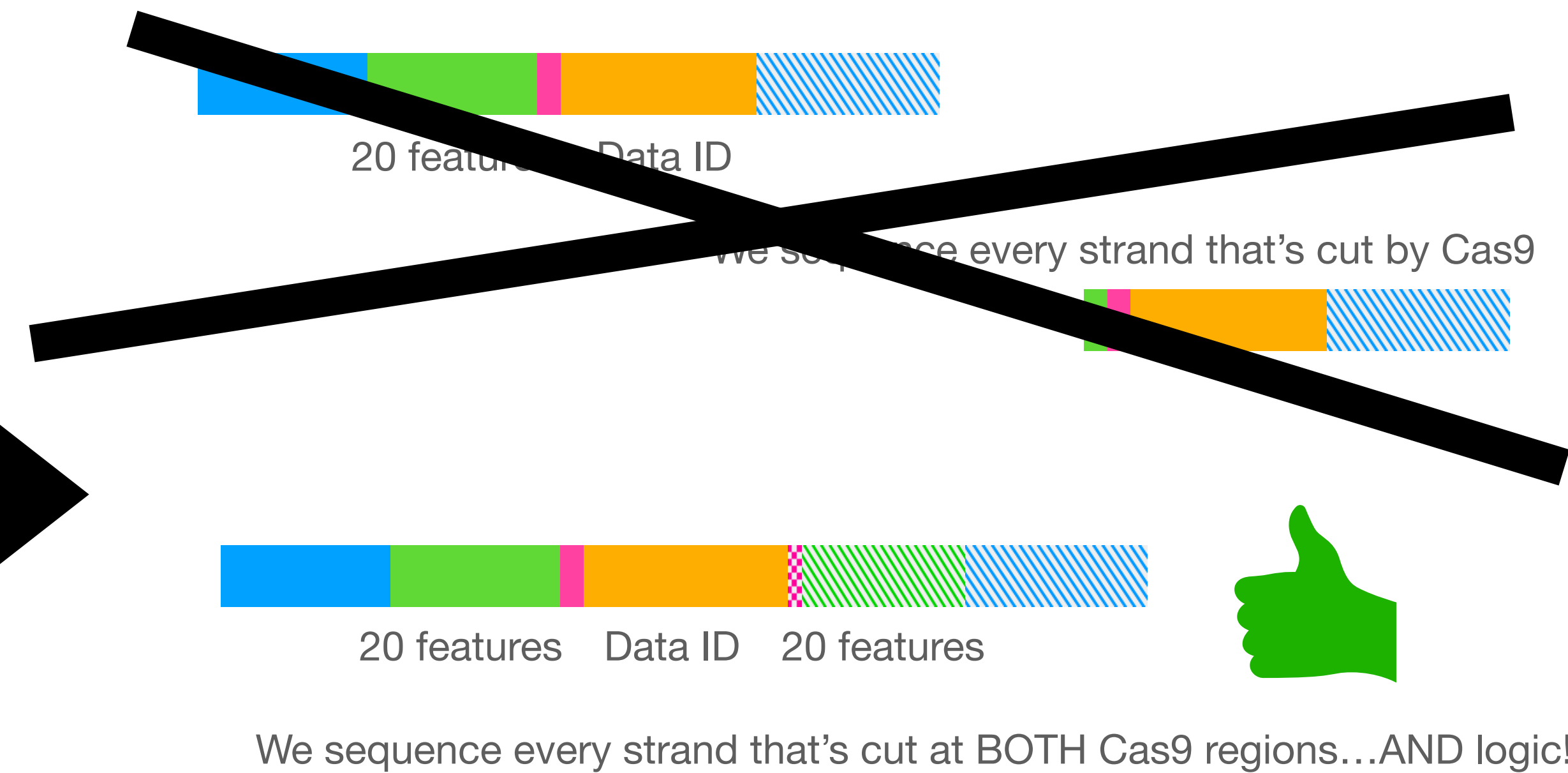
Using Cas9

C. Training loop



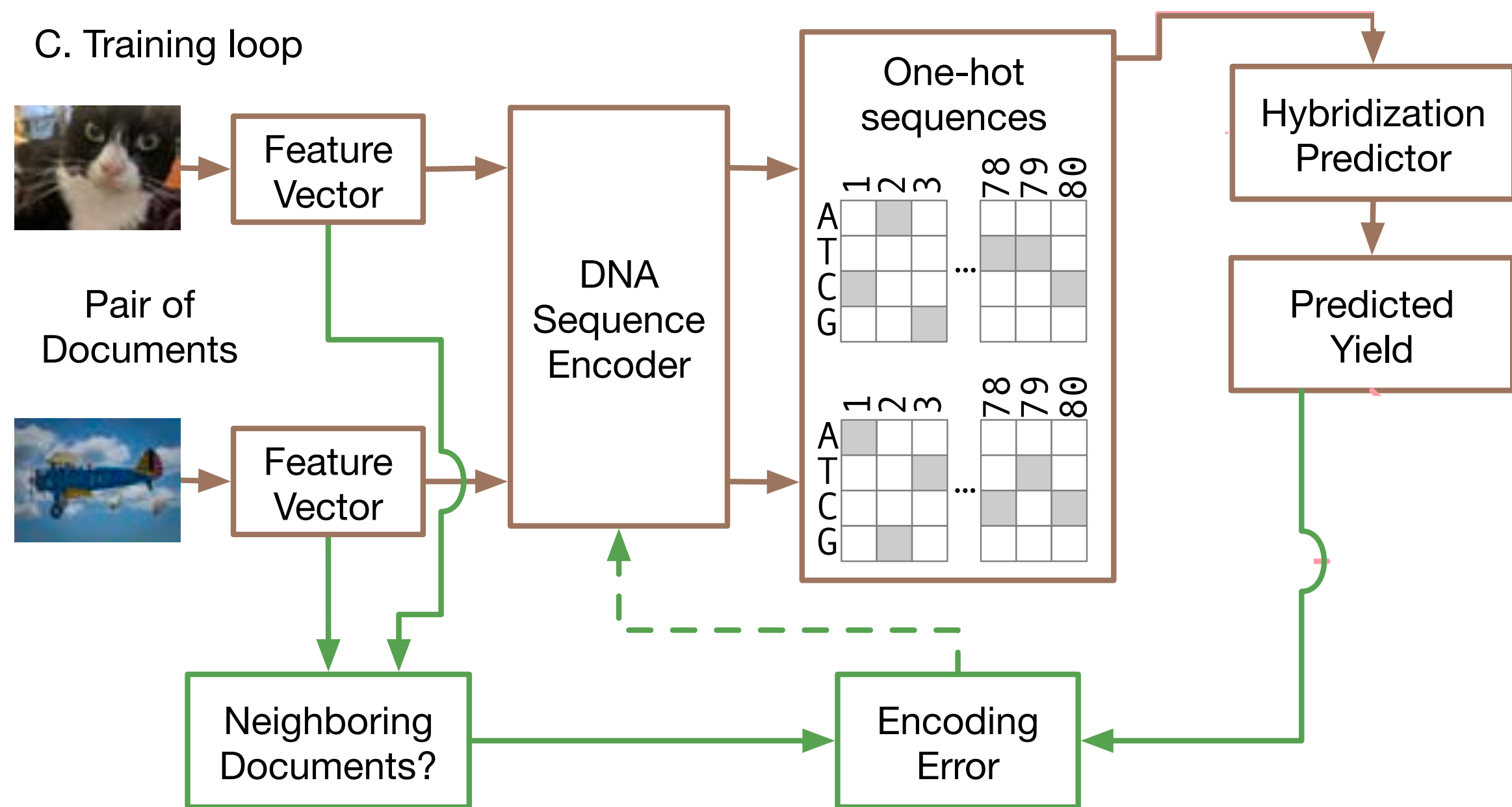
Cas9

Hard to do in wet lab



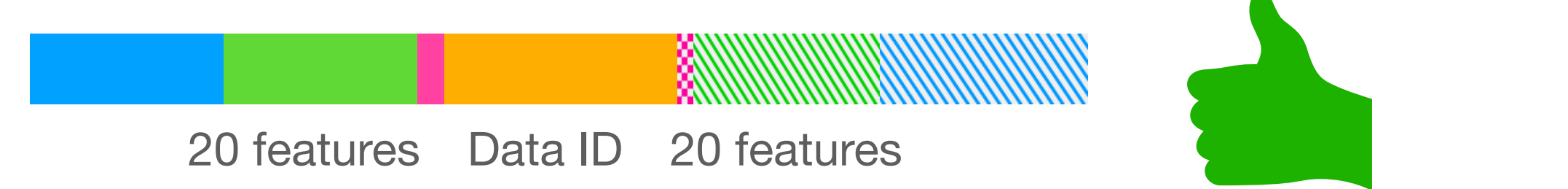
Using Cas9

C. Training loop



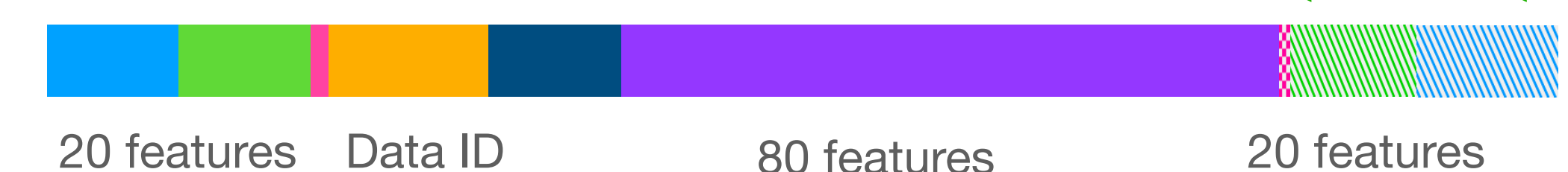
Cas9

Hard to do in wet lab



We sequence every strand that's cut at BOTH Cas9 regions...AND logic!

What if we allow for nested queries??

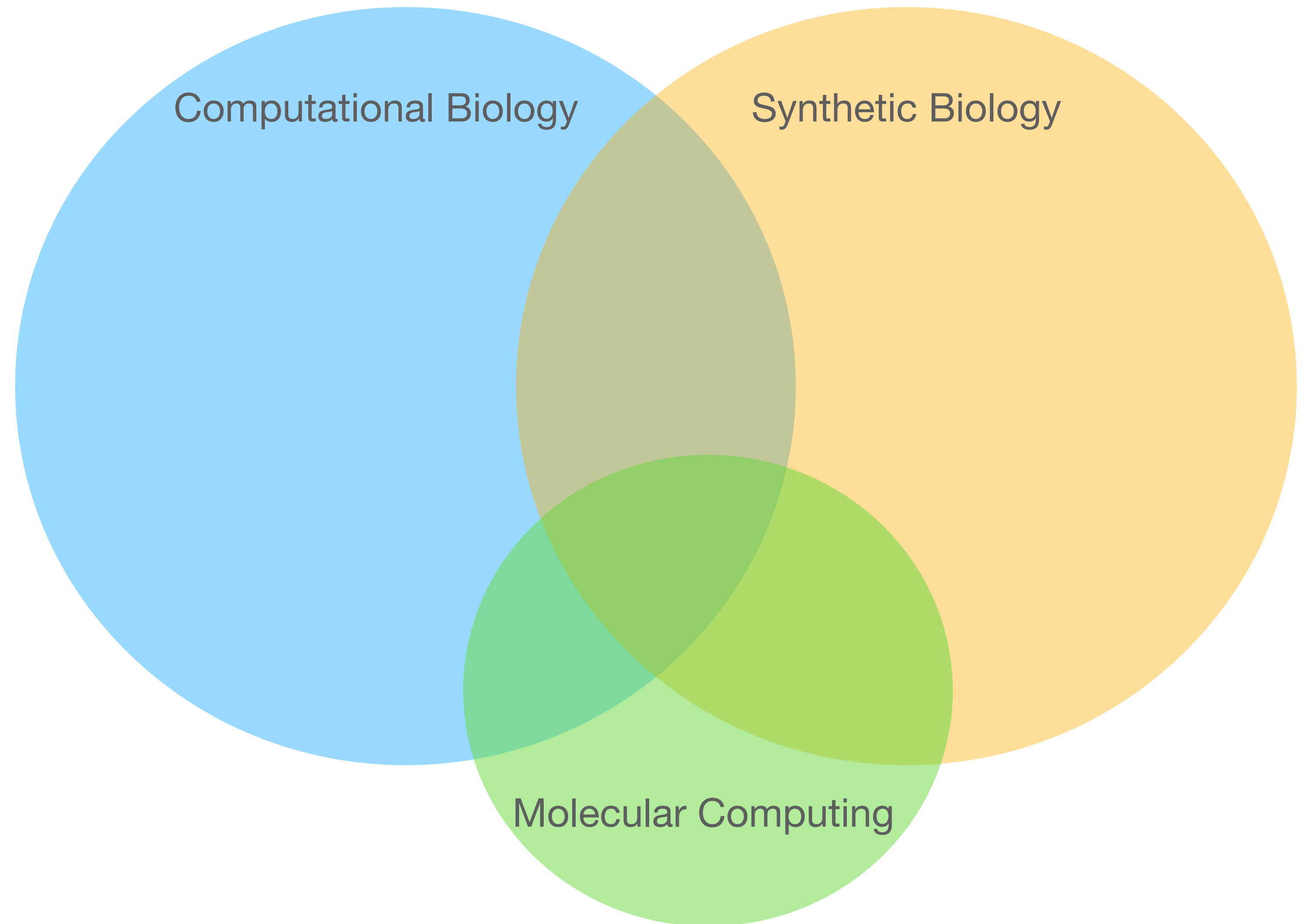


Comp Bio \cap Molecular Computing \cap SynBio

The lines are still being drawn

My two cents for people in any of these three fields:

- If you're a computation-centered person, get comfortable talking to wet lab-centered people
- vice-versa
- Bigger computational AND molecular toolboxes tend to make it easier to design experiments
- We need more tools



Comp Bio ∩

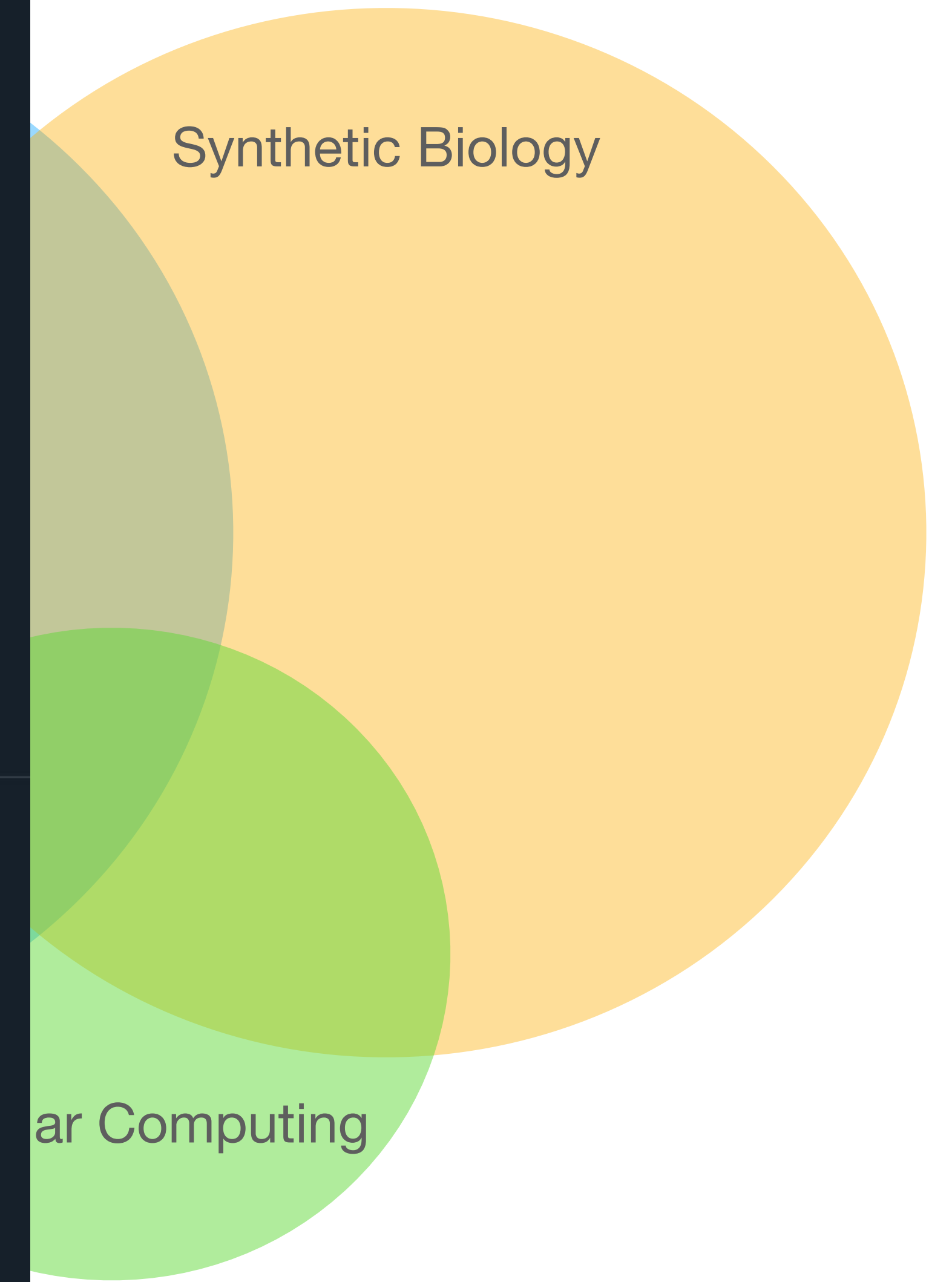
The lines are still being drawn

My two cents for people in a
these three fields:

- If you're a computation-c person, get comfortable t wet lab-centered people
- vice-versa
- Bigger computational AN molecular toolboxes tend make it easier to design experiments
- We need more tools

The screenshot shows a Twitter thread on a dark background. At the top, the word "Tweet" is centered. The first tweet is from Patrick Boyle (@p_maverick_b) with the text "Cries in Synthetic Biologist". Below it is a reply from Kyle (@KyleMorgenstein) with the text "how the hell did we make planes before CAD??". The tweet shows 5 retweets and 31 likes. Below that is a reply from Sebastian S. Cocioba (@ATinyGre...) asking "What's the slide-rule of synthetic biology tho?". The final tweet is from Patrick Boyle (@p_maverick_b) replying to Sebastian, with the text "Back in my day we tagged everything with GFP whether you needed it to glow or not".

ing ∩ SynBio



Things to think about

- What are some tools you'd like to see developed?
- What are things you'd like to see standardized?
- Are there times when having a deeper background in a different field (i.e., biology) would have helped you?
- Anything other wishes for the future?