

^eMicrosoft Research

Generation of Sequence–Activity Datasets for SlugCas9 Assays Using Diverse PAMs Using Multiplexed-Barcoded Sequence Display

A Recording plasmid

SlugCas9

P21199

Barcode DNA

Targeting plasmid

sgRNA

Target DNA

Recording Barcode

SlugCas9 Variants

PAM

NNGA

NNGT

NNGC

NNGG

Next-Generation Sequencing

Average Mutation Number in Barcodes

Sequence	NNGA	NNGT	NNGC	NNGG
SlugCas9-1	3.3	1.7	1.8	0.9
SlugCas9-2	1.8	0.8	2.1	3.7
SlugCas9-3	2.9	2.8	3.1	1.9

Sequence-Activity Data for SlugCas9

Machine Learning

Hit with General PAM Activity

Hit towards NNGA PAM

Prediction

Ground Truth

NNGA

NNGT

NNGC

NNGG

B Docking Model

Interface - Zoom In View

Binding Pocket - Side View

SlugCas9

Barcode DNA & NNGG-PAM

M990

K1016

E1012

S985

N984

C Reporting Plasmid

P21199

P21

2 Stop Codons

sgGFP

Editing Plasmid

P2190

TadA8e

SlugCas9 Variant

Translation

sfGFP Protein

4321 Lys Met Phe Glu Stop Stop Leu

5'-NCNN AACATG TTT GAG TCA TCA CTT -3'

3'-NNGN TTN TAC AAA CTG ACT TGT GAA -5'

Target

TadA8e

SlugCas9

Base Editing

A to G Mutation

4321 Lys Met Phe Glu Arg Gln Leu

5'-NCNN AACATG TTT GAG CGC GAA CTT -3'

3'-NNGN TTN TAC AAA CTA CTG CT TGT GAA -5'

D Average Mutation Number

NNGA NNGT NNGC NNGG

SlugCas9-WT

SlugCas9-N984S

SlugCas9-K1016I

SlugCas9-Dead Variant

E Fluorescence/OD₆₀₀

NNGA NNGT NNGC NNGG

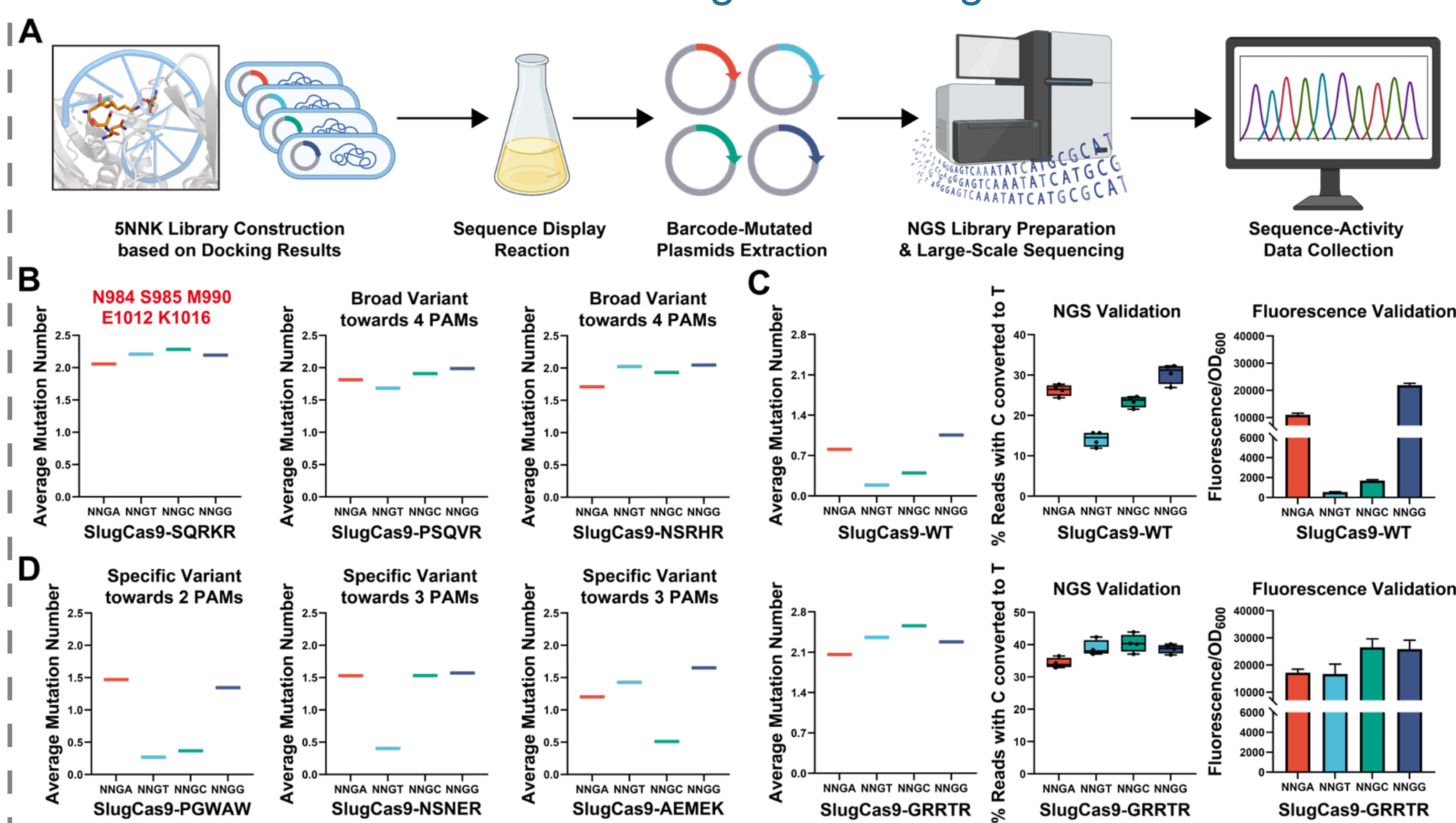
SlugCas9-WT

SlugCas9-N984S

SlugCas9-K1016I

SlugCas9-Dead Variant

Generating Large-scale sequence-activity datasets for SlugCas9 to enable machine learning-based SlugCas9 evolution



JOHN S. DUNN
FOUNDATION

Engineering proteins with desired functions remains challenging due to the labor-intensive nature of traditional methods and limited sequence–activity data for machine learning. Here, we present Sequence Display, a scalable platform that rapidly generates large-scale protein sequence–activity datasets, enabling machine learning-driven protein evolution. Sequence Display can be multiplexed to assess the specificity of individual mutants within a single experiment. By integrating these datasets with pre-trained protein language models, we construct fine-grained, variant-specific activity landscapes to identify high-performance variants. We demonstrate its broad applicability by generating datasets for cytosine deaminase, uracil glycosylase inhibitor, and a compact Cas9 nuclease. For Cas9, we produced over 32 million data points and evolved a variant with expanded PAM recognition, which outperformed a previously reported mutant from phage-assisted evolution. This study establishes Sequence Display as a powerful tool for mapping sequence–function relationships and accelerating the discovery of optimized proteins for biological and medical applications.

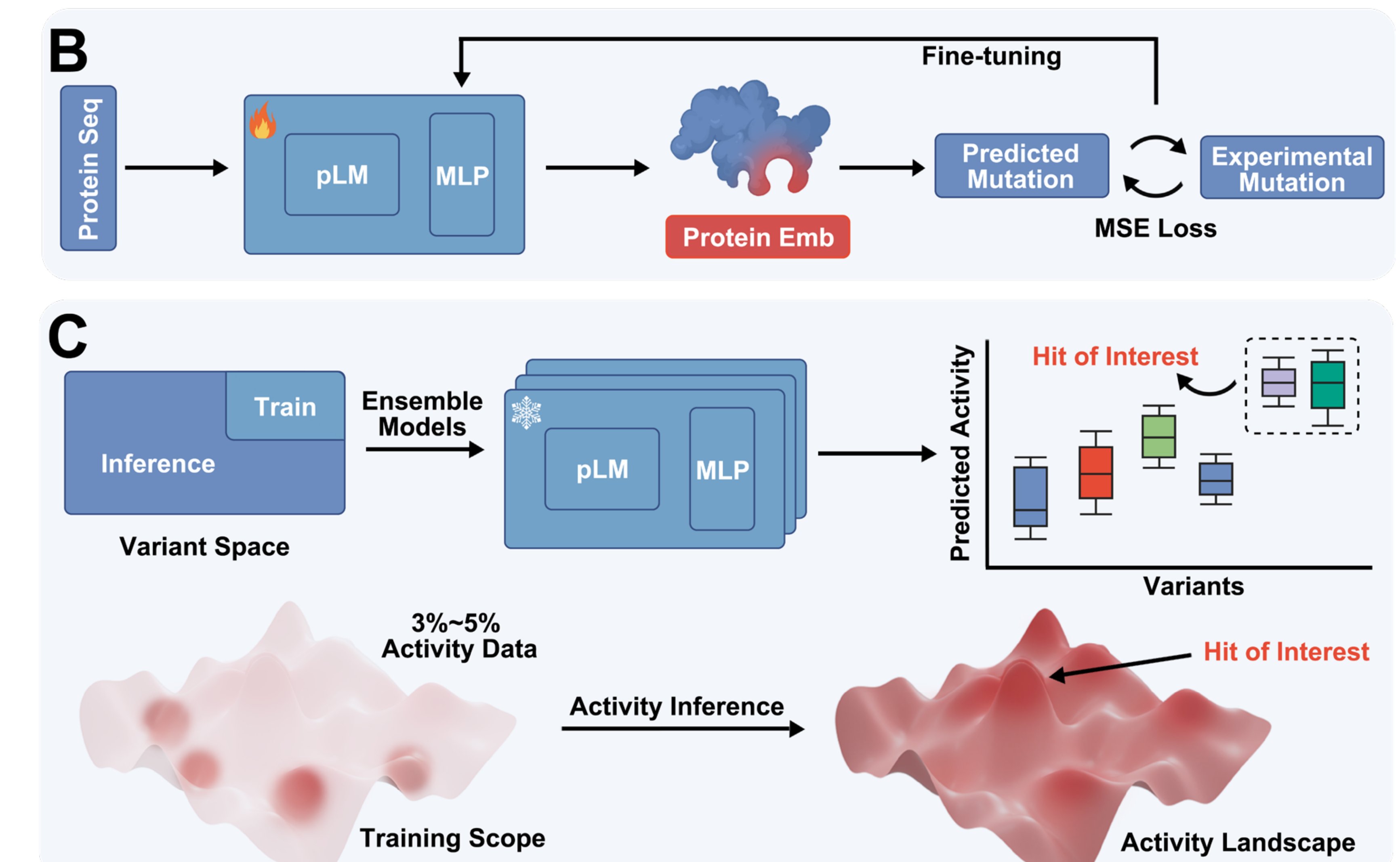
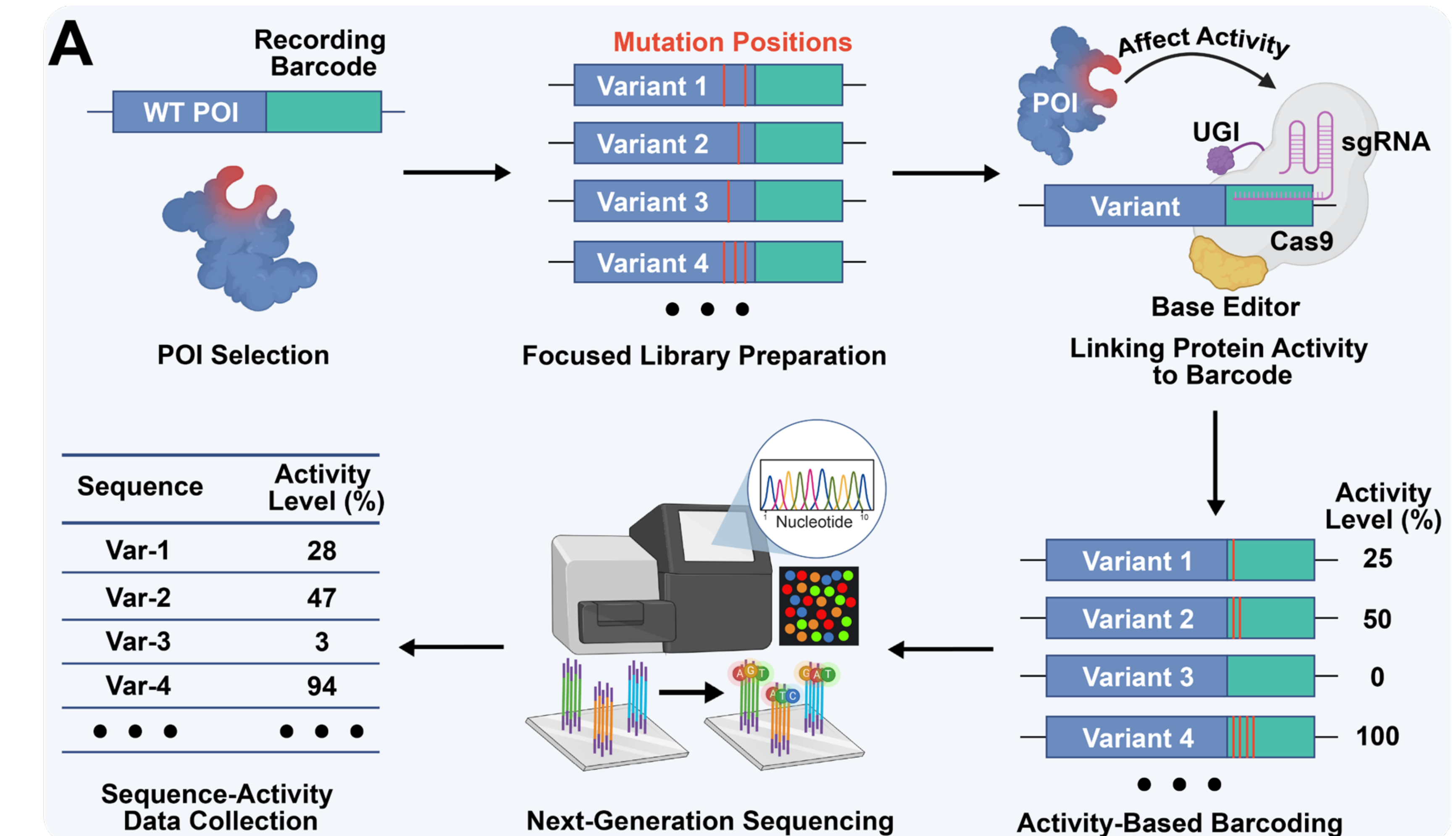


Fig. 1. Schematic overview of Sequence Display-enabled machine learning pipeline. (A) Sequence Display platform. (B) Machine learning framework for Sequence Display protein evolution. (C) Ensemble-based activity inference and landscape construction.

Modeling the Relationship between SlugCas9 Variants and their Activities using pLMs

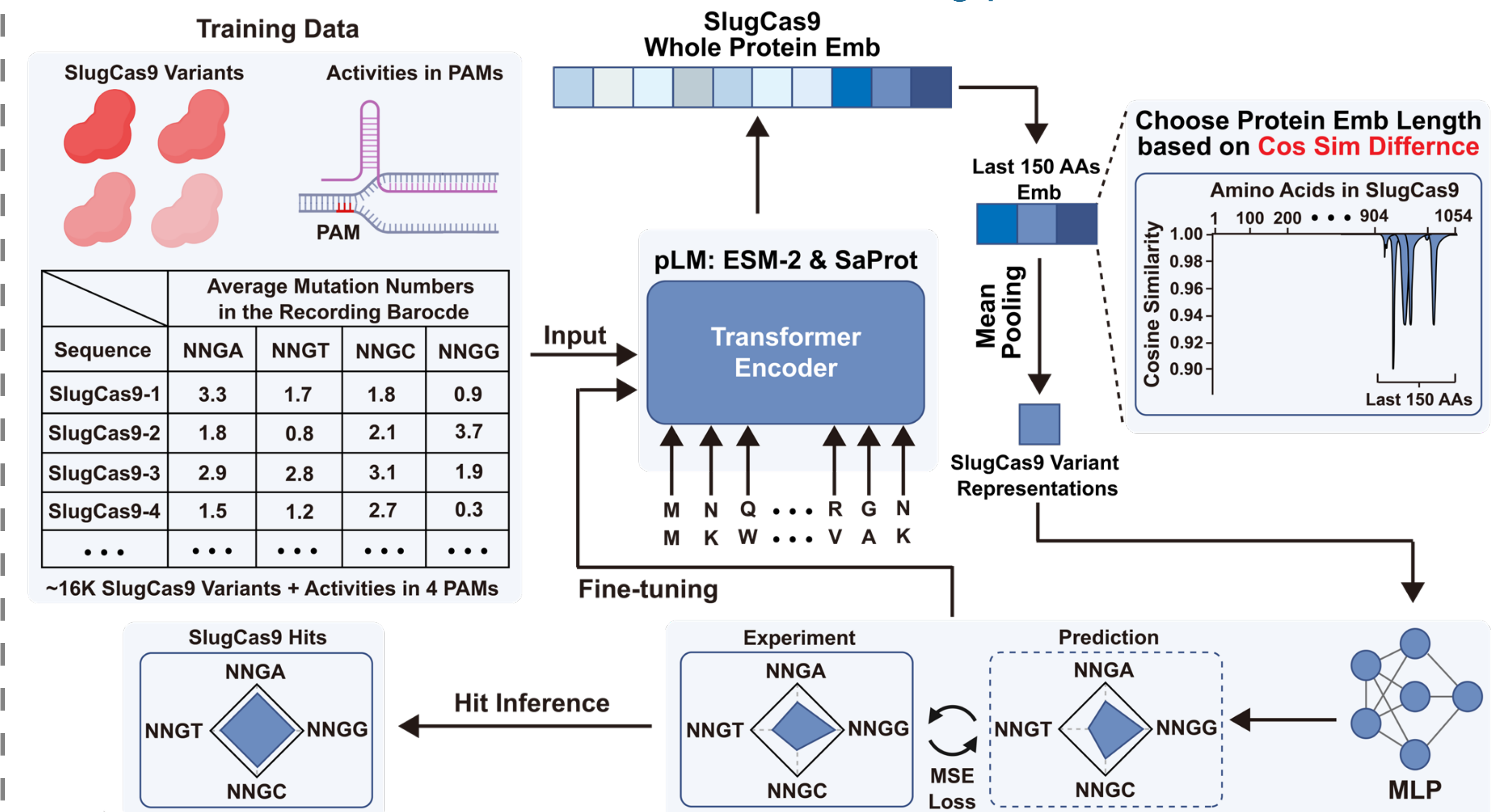


Fig. 4. Modeling the activity landscape of SlugCas9 with protein language models. Machine learning pipeline for predicting SlugCas9 activity. pLMs, including ESM-2 and SaProt, are used to extract local embeddings based on cosine similarity calculations for each amino acid in the SlugCas9 sequence.

Fig. 5. Validation of ensemble models-predicted SlugCas9 hits and construction of the activity landscape. (A) High-throughput validation for top-predicted SlugCas9 variants. (B) High-throughput validation of an additional batch of top-predicted variants. (C) NGS-based validation of top-performing variants. (D) Fluorescence-based validation of top-performing variants. (E) UMAP visualization of the SlugCas9 variant activity landscape.