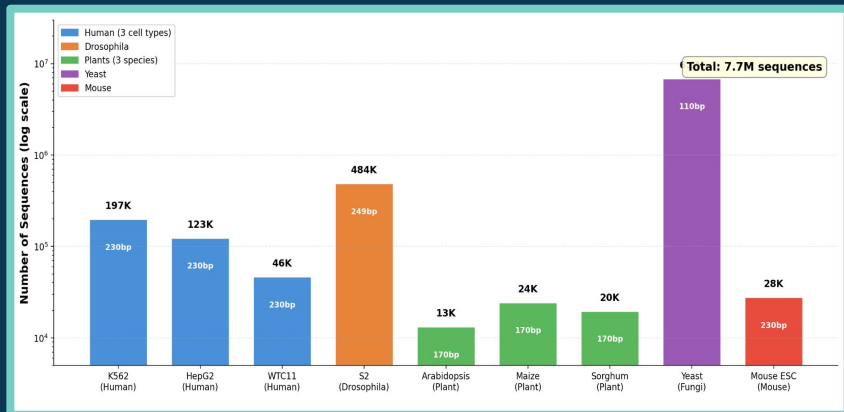# Sequence Data



**Figure 22,23,24.** Dataset summary of 7.5+ million sequences across seven species: human cell lines (K562, HepG2, WTC11)(ENCODE Project Consortium, 2020; Inoue et al., 2017), Drosophila S2(de Almeida et al., 2022; Arnold et al., 2013), plants (Arabidopsis, Maize, Sorghum)(Jores et al., 2021), yeast(de Boer et al., 2020; Schreiber et al., 2020), and mouse ESC(Kalkan et al., 2017), with sequence counts per dataset. **(Figure generated by student author using python, matplotlib)**
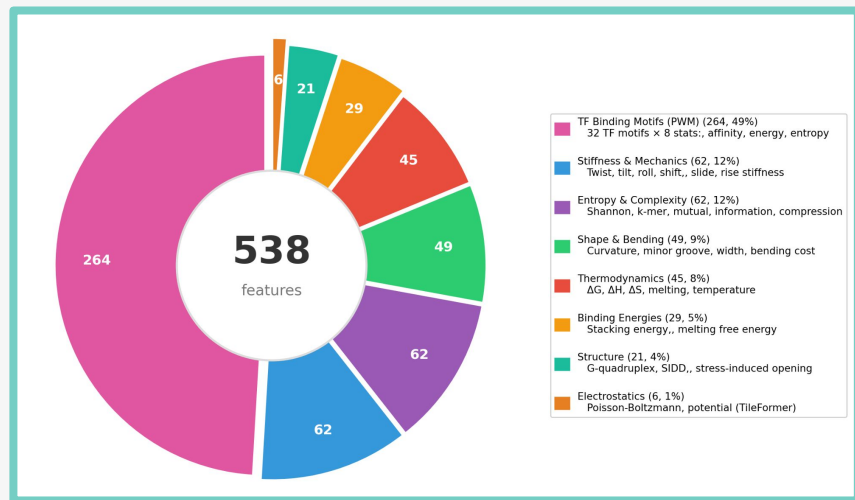
# Biophysics Training Data



**Figure 25**. Feature breakdown with TF features and pure biophysical features **(Figure generated by student author using python, matplotlib)**
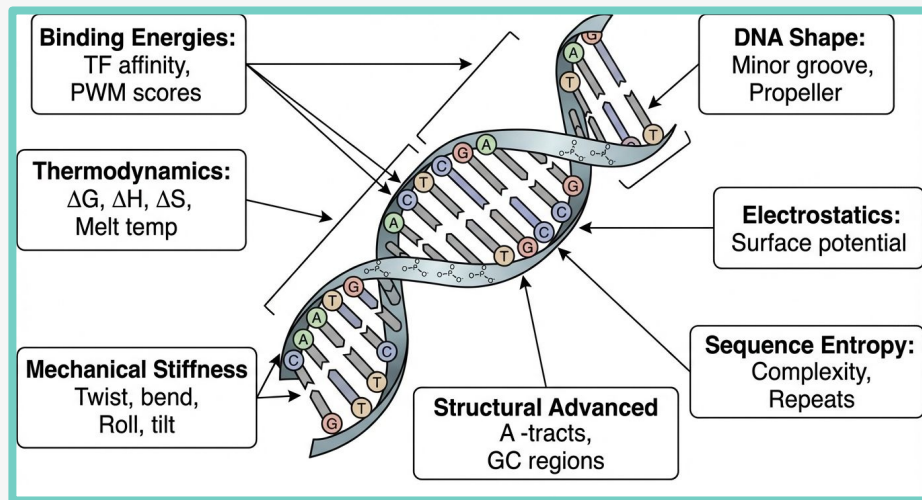


**Figure 26**. Representation of DNA and where the features are found physically **(Figure generated by student author using FigureLabs)**
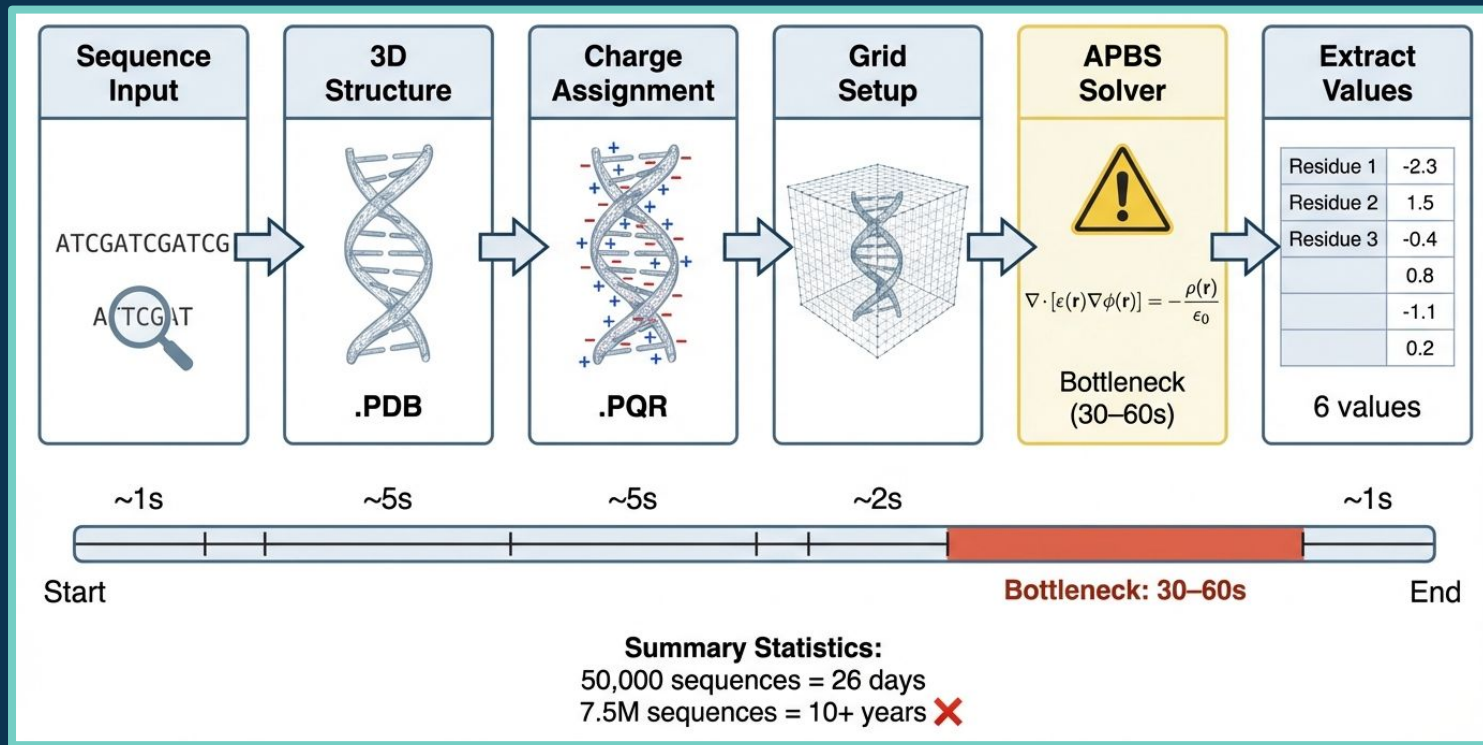
# Problem: Manual Labelling



| Sequence Input | 3D Structure | Charge Assignment | Grid Setup | APBS Solver | Extract Values |
|---|---|---|---|---|---|

ATCGATCGATCG

ATCGAT

.PDB

.PQR

$$\nabla \cdot [\varepsilon(\mathbf{r})\nabla\phi(\mathbf{r})] = -\frac{\rho(\mathbf{r})}{\varepsilon_0}$$

Bottleneck (30–60s)

| Residue 1 | -2.3 |
| Residue 2 | 1.5 |
| Residue 3 | -0.4 |
| | 0.8 |
| | -1.1 |
| | 0.2 |

6 values

~1s      ~5s      ~5s      ~2s      ~1s

Start

Bottleneck: 30–60s

End

**Summary Statistics:**
50,000 sequences = 26 days
7.5M sequences = 10+ years ❌

**Figure 27:** Manually labelling with solving Poisson–Boltzmann is time–consuming**(Figure generated by student author using FigureLabs)**
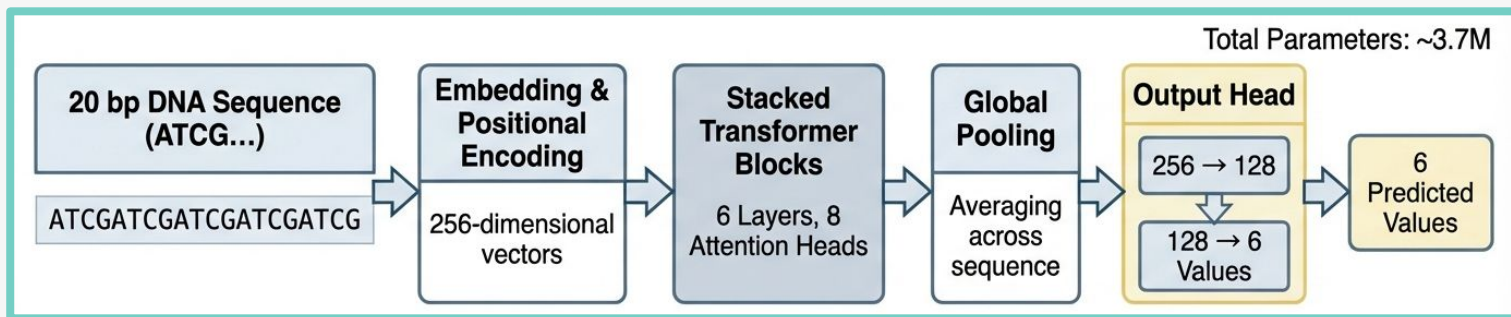
# Solution: TileFormer



**Figure 28.** TileFormer architecture diagram from input to predicted values**(Figure generated by student author using FigureLabs)**



**Figure 29.** TileFormer output breakdown**(Figure generated by student author using FigureLabs)**
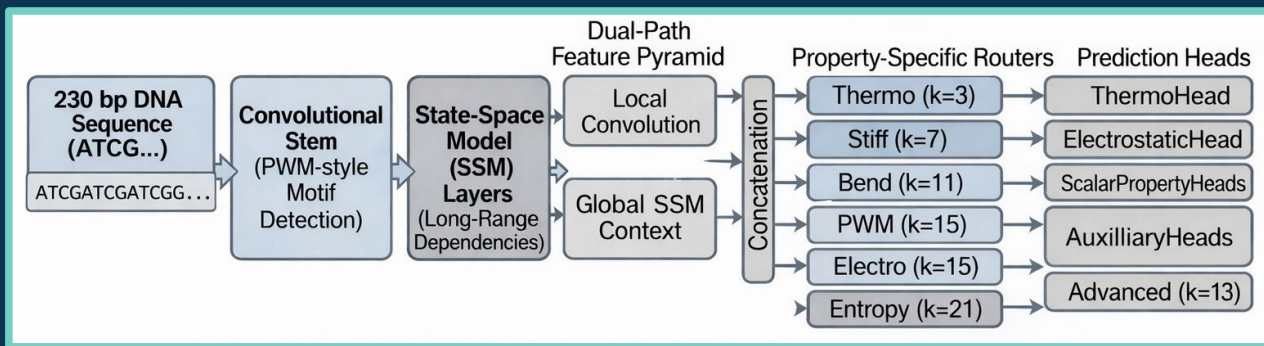
# Solution 2: PhysInformer



**Figure 30.** PhysInformer architecture diagram from input to predicted values**(Figure generated by student author using FigureLabs)**
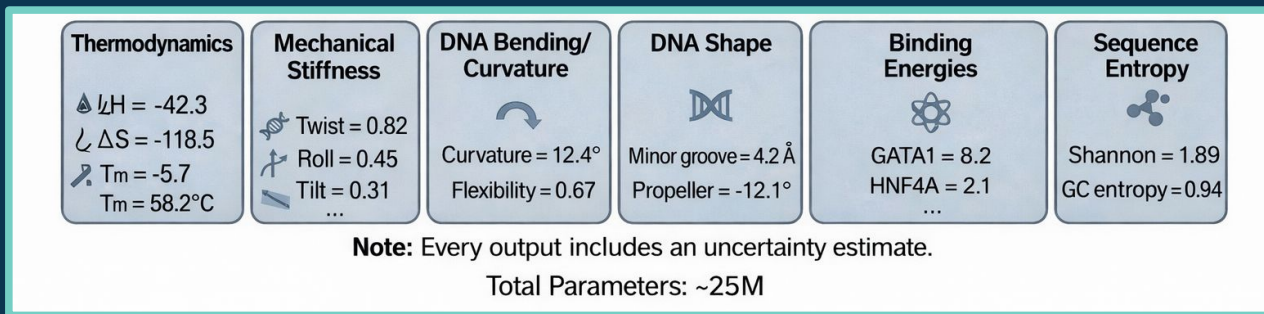


**Figure 31.** PhysInformer output breakdown**(Figure generated by student author using FigureLabs)**

# CADENCE: **C**onvolutional **A**rchitecture for **D**NA **EN**hancer **C**haracterization and **E**xpression prediction
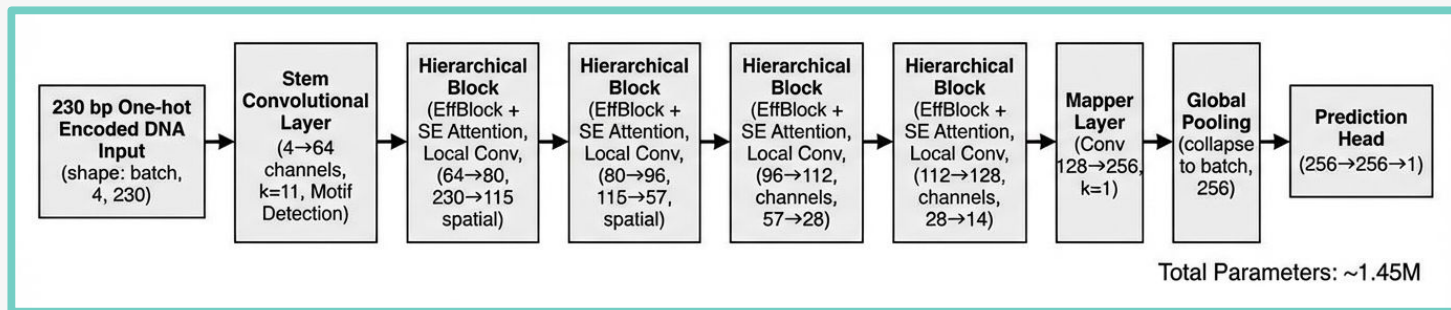


**Figure 32.** CADENCE architecture diagram from input to predicted values**(Figure generated by student author using FigureLabs)**
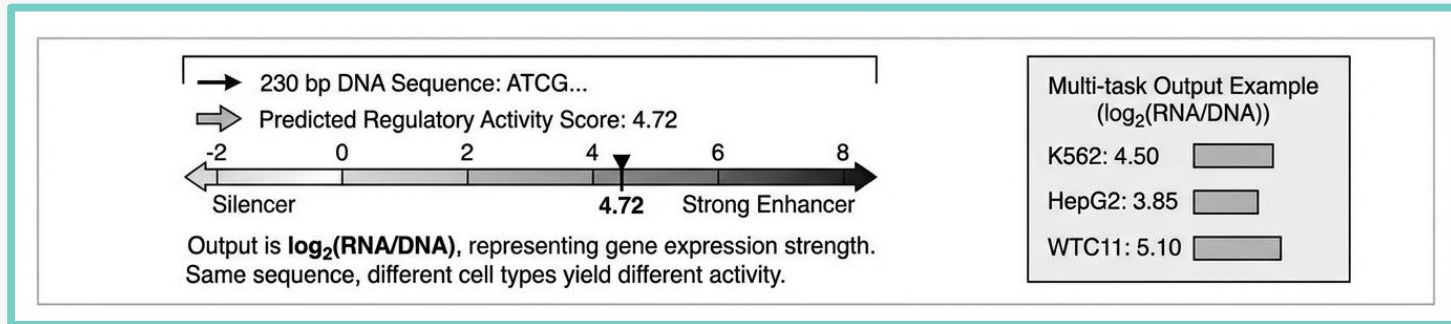


**Figure 33.** CADENCE output breakdown**(Figure generated by student author using FigureLabs)**

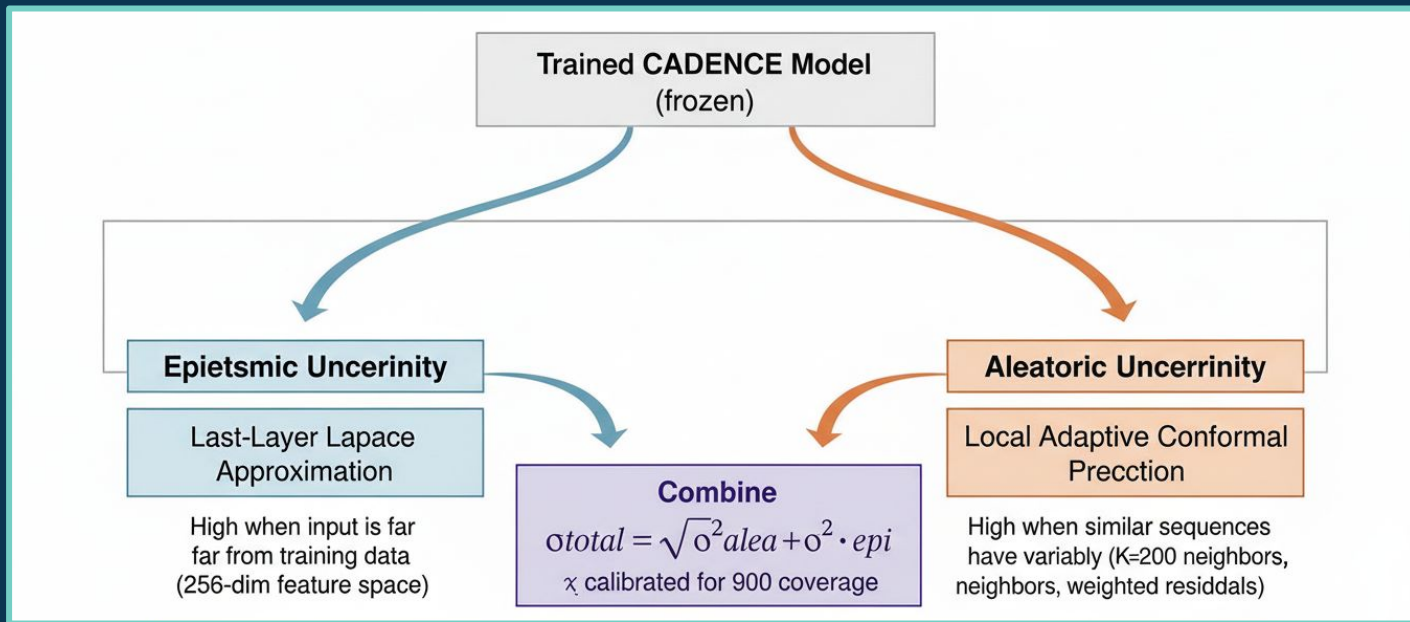# PLACE: Post-hoc Laplace And Conformal Estimation



**Figure 34.** PLACE framework diagram: frozen CADENCE model splits into two parallel paths—epistemic uncertainty via last–layer Laplace approximation ($\Sigma = (\lambda I + \sigma^{-2}\Phi^T\Phi)^{-1}$) and aleatoric uncertainty via K–NN conformal prediction—combined as $\sigma\_total = \sqrt{(\sigma^2\_alea + \kappa \cdot \sigma^2\_epi)}$ with $\kappa$ calibrated for 90% coverage. **(Figure generated by student author using FigureLabs)**
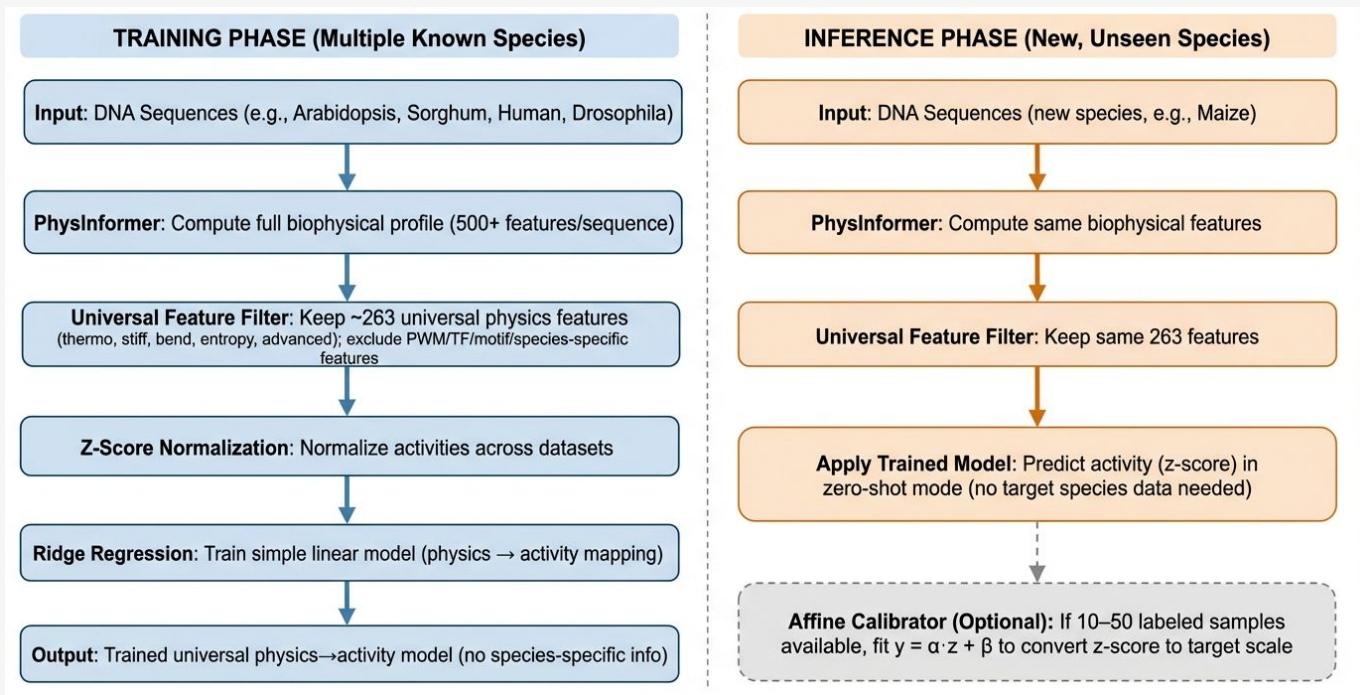
# S2A: Sequence to Activity



**Figure 35.** S2A pipeline flowchart: training species sequences → PhysInformer → universal feature filter (keep ~263 physics features, exclude PWM/TF binding) → z-score normalization → Ridge regression → trained model applied zero-shot to held-out species → activity prediction without target-species data. **(Figure generated by student author using FigureLabs)**
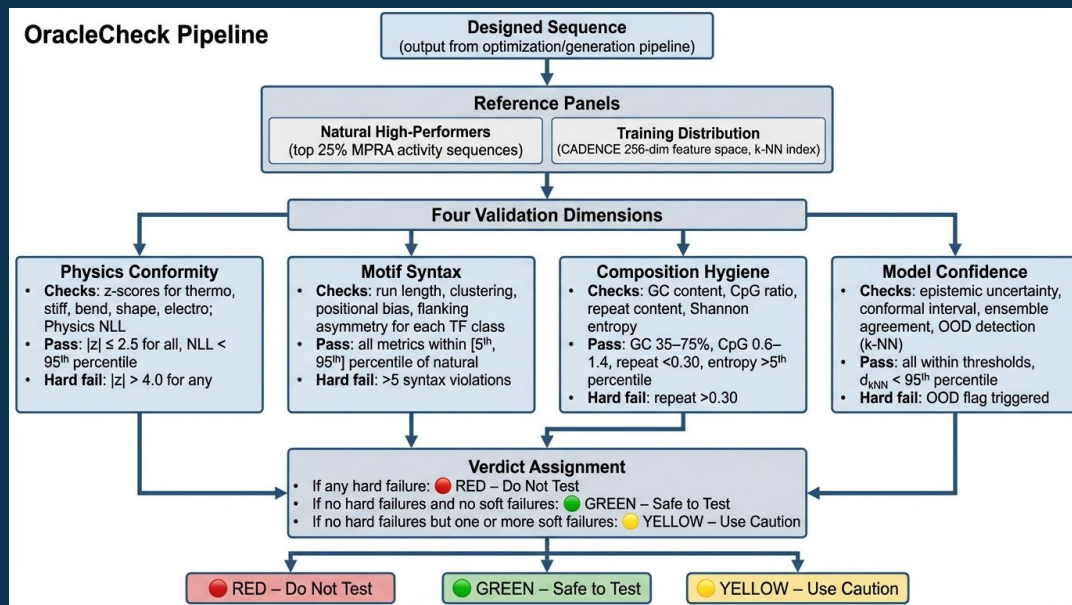
# OracleCheck



**Figure 36.** OracleCheck validation flowchart: designed sequence evaluated across four dimensions—(1) Physics Conformity ($|z| \leq 2.5$), (2) Motif Syntax ([5th, 95th] percentile), (3) Composition Hygiene (35% ≤ GC ≤ 75%, repeat < 30%), (4) Model Confidence ($\sigma\_epi$ < 90th %ile, OOD check)—yielding GREEN (pass all), YELLOW (1 soft fail), or RED (hard fail) verdict. **(Figure generated by student author using FigureLabs)**
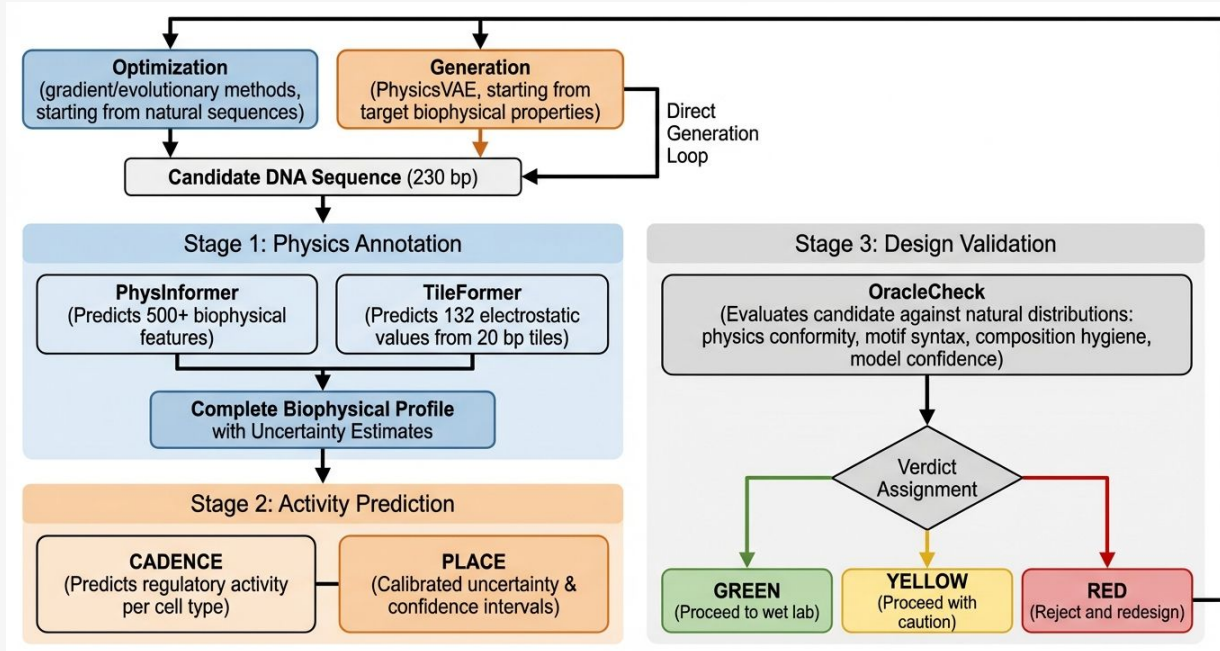
# Integrated Pipeline



**Figure 37.** FUSEMAP integrated loop diagram: candidate sequence → PhysInformer + TileFormer (physics annotation) → CADENCE + PLACE (activity prediction with uncertainty) → OracleCheck (naturality validation) → GREEN sequences proceed to experimental validation; RED sequences loop back to redesign. PhysicsVAE shown as alternative entry point for physics–conditioned sequence generation. **(Figure generated by student author using FigureLabs)**