

Predicting DNA function when grafted across species

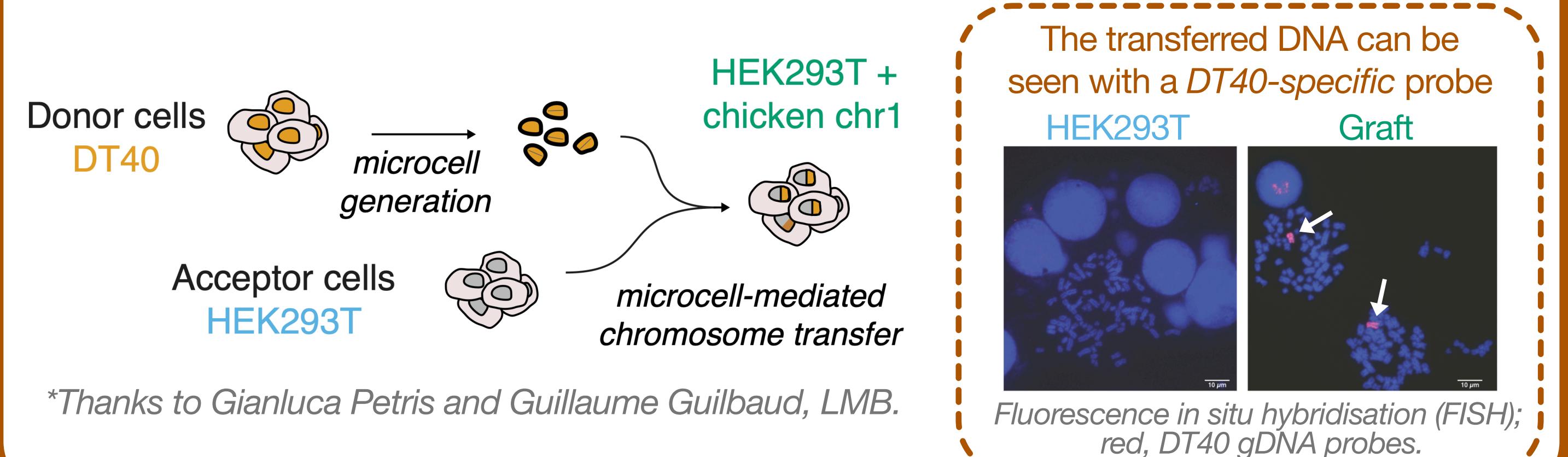
Sergio Garcia Busto¹, Jacob Hepkema¹, Pierre Murat¹, Leopold Parts^{1,2}

¹Wellcome Sanger Institute, Wellcome Genome Campus, Hinxton, Cambridgeshire, CB10 1SA

²ExpressionEdits, Haverhill Research Park, Haverhill, CB9 7LR

① The chicken graft: a foreign chromosome

We have **human HEK293T cells (hg38)** with a **chromosome from a chicken cell line, DT40 (gg6)**. This was generated via microcell-mediated chromosome transfer by our collaborators* at the Laboratory of Molecular Biology (LMB). The foreign chromosome (gg6_1) is **197 Mb long**, and has been characterised extensively in the Parts lab.



Fluorescence in situ hybridisation (FISH); red, DT40 gDNA probes.

② Why should we try to predict graft data?

Deep learning (DL) models trained on genomics attempt to learn a sequence function or DNA 'grammar'. But the chicken graft has **two distinct characteristics**:

Sequence divergence through 319 million years

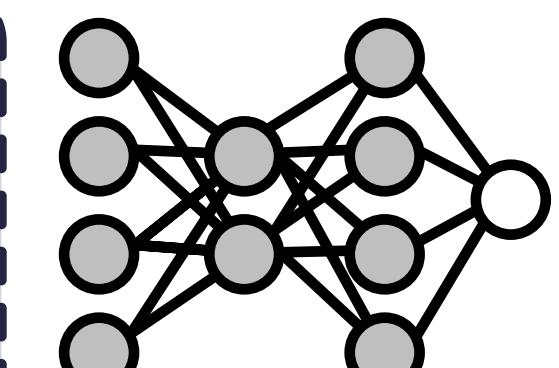
Foreign cellular context

By evaluating how DL models predict graft data when exposed to different sequences, we can **assess the extent to which grafted grammar is learnable**

This poster focuses on **H3K9me3**, a histone mark associated with **heterochromatin**.

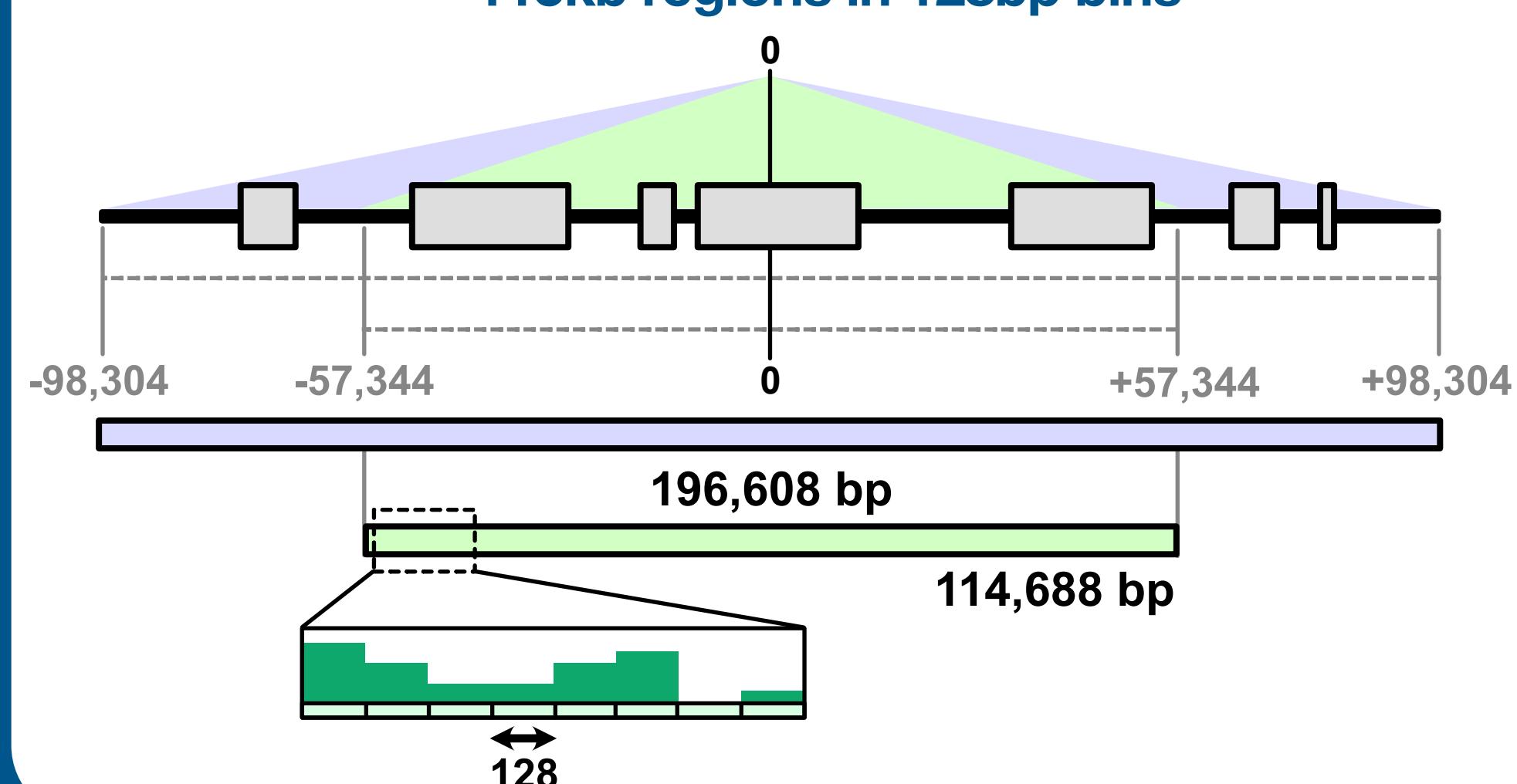
Aims:

- (i) Evaluate the ability of DL models to **learn regulatory grammar of the grafted chromosome** using H3K9me3
- (ii) Interpret results by fine-tuning different models and **comparing predictions to data from native and grafted contexts**

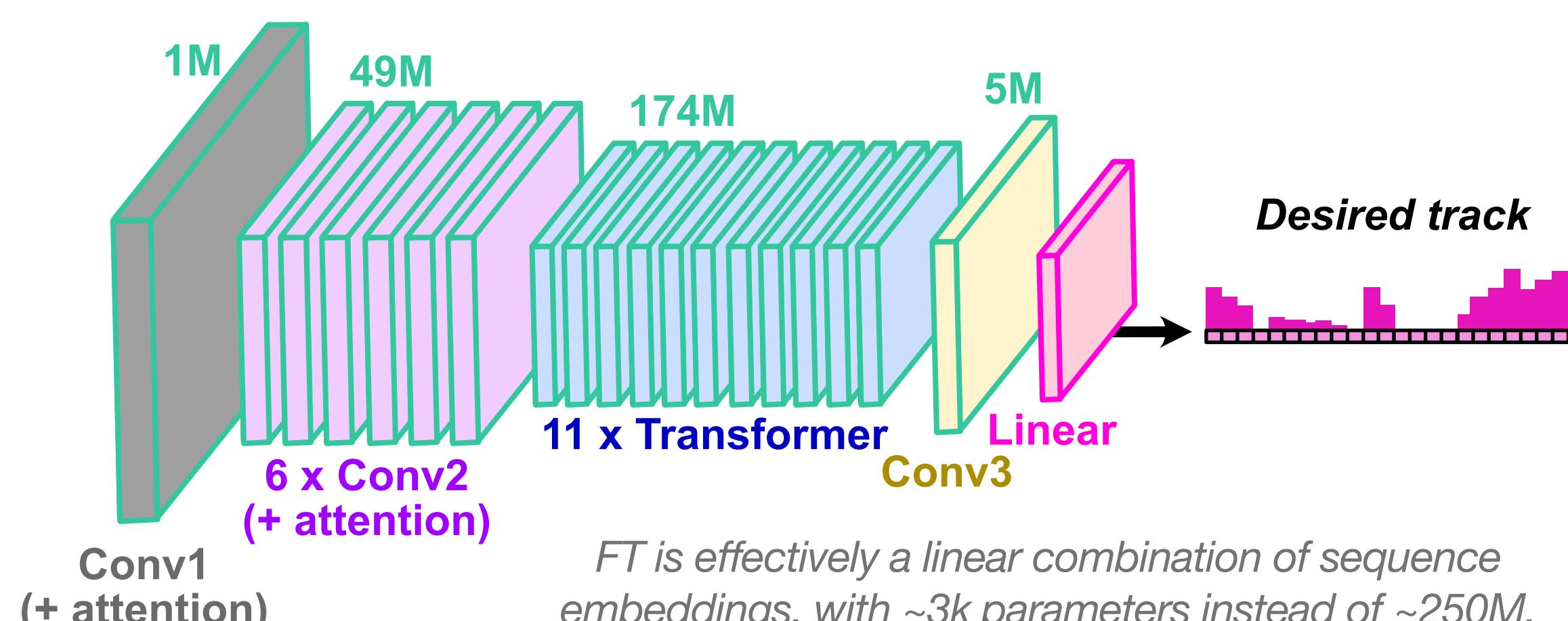


③ Enformer: a Large Language Model (LLM) that predicts genomic tracks

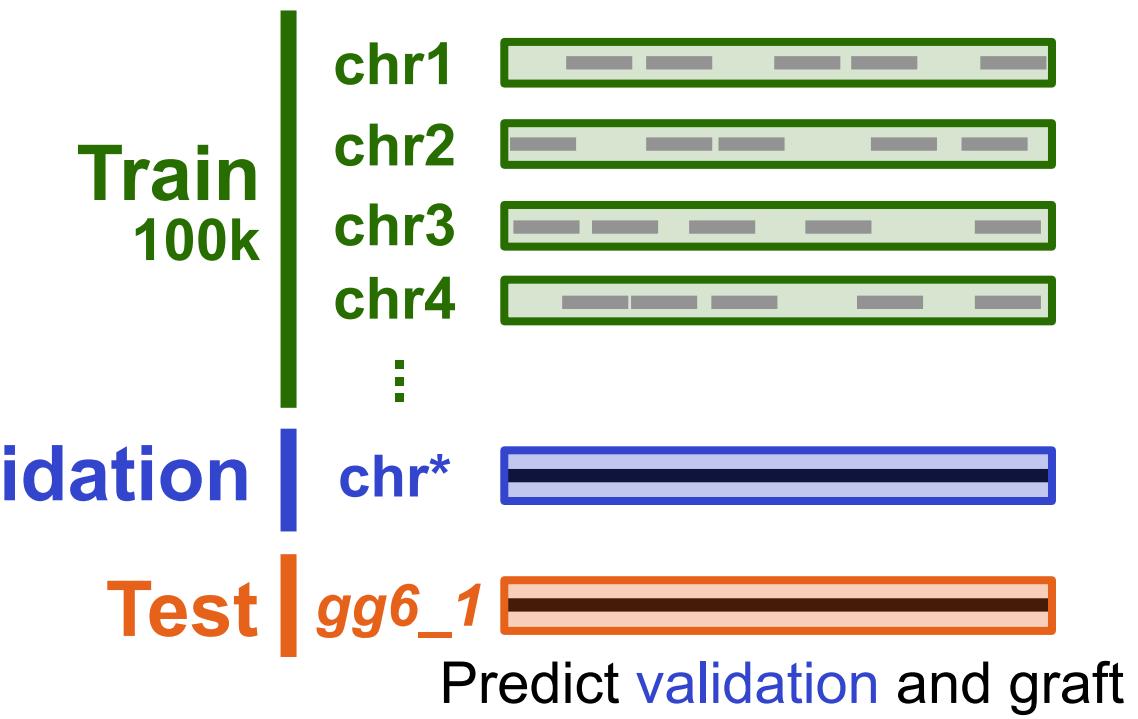
Enformer is trained on **human and mouse genomic tracks**, has a **context window of ~200kb** and predicts the central **~115kb regions in 128bp bins**



Enformer **fine-tuning (FT)**: model **weights are frozen**, and a **custom linear layer is added and trained**. Learned sequence features (embeddings) are used to learn desired grammar.

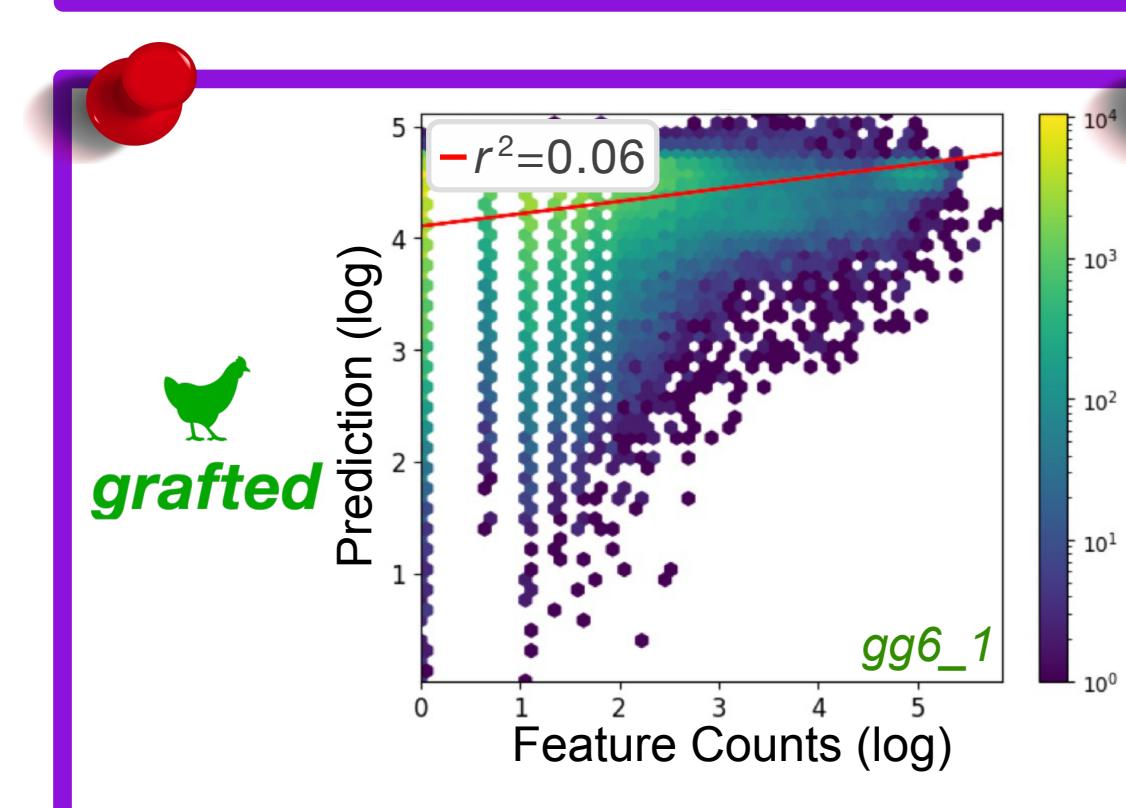
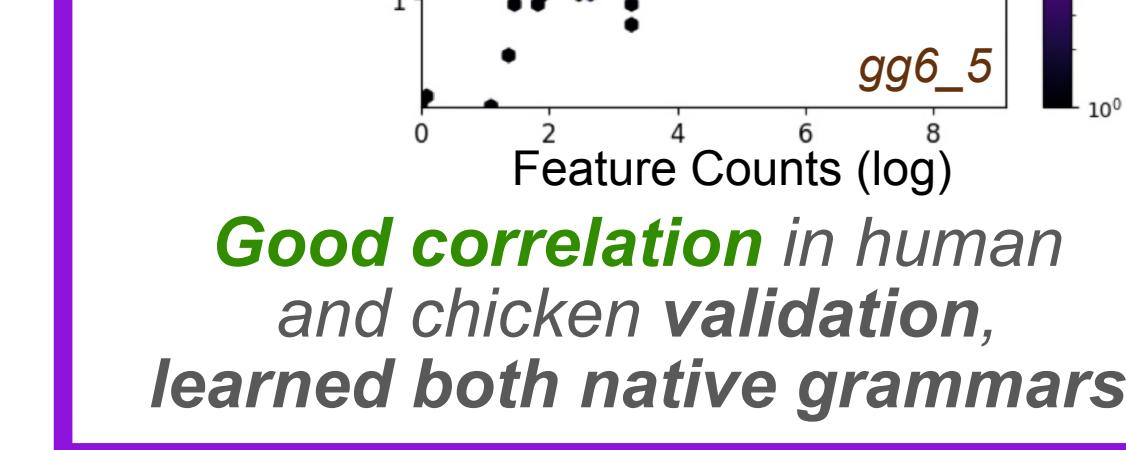
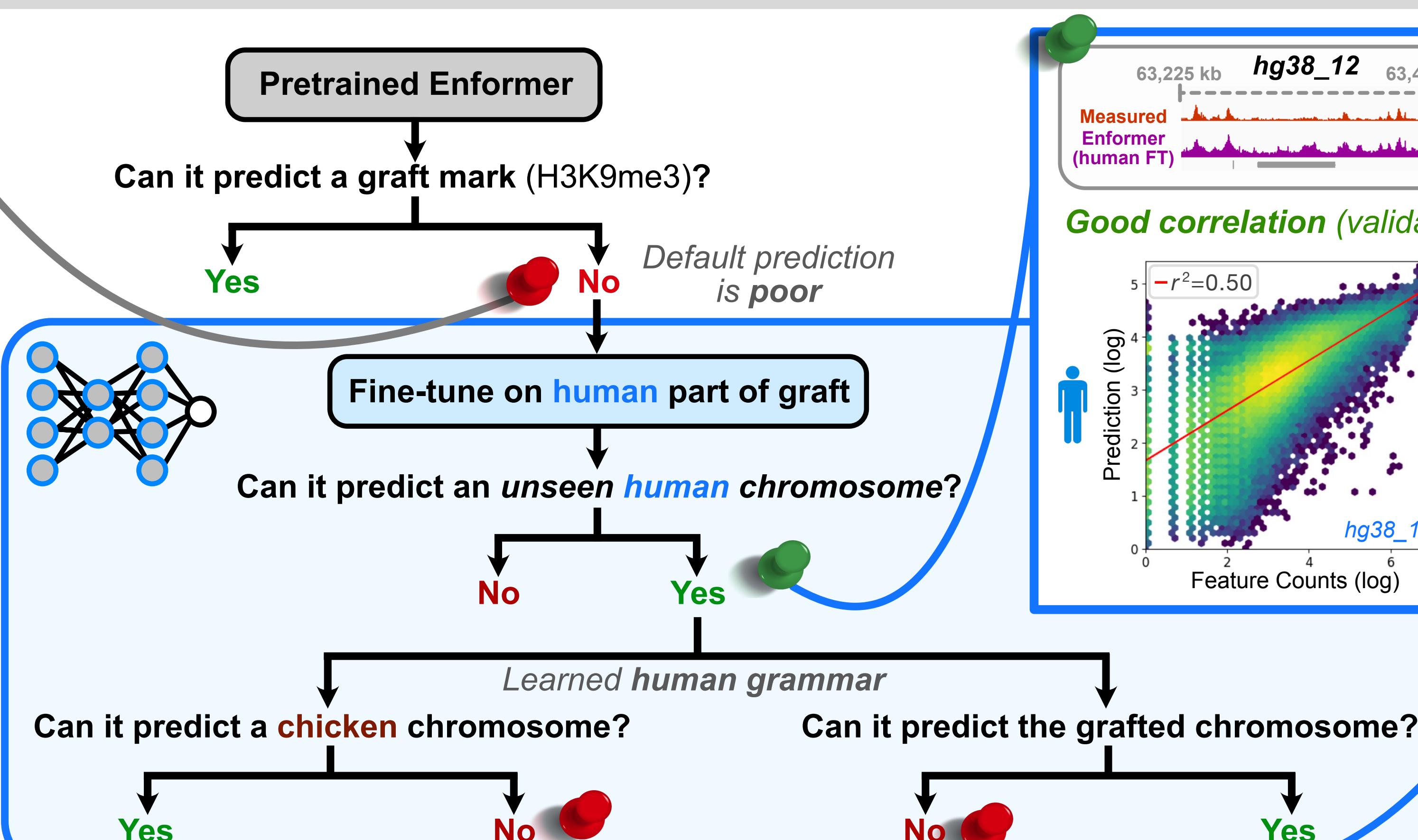
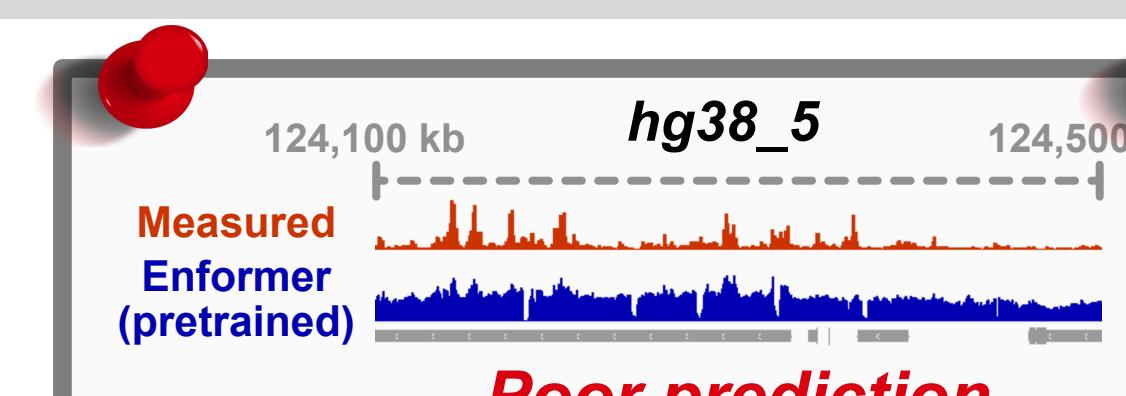


For **training**, 100,000 sequences of ~200kb are randomly sampled, and an **entire chromosome is held out** to evaluate FT efficacy

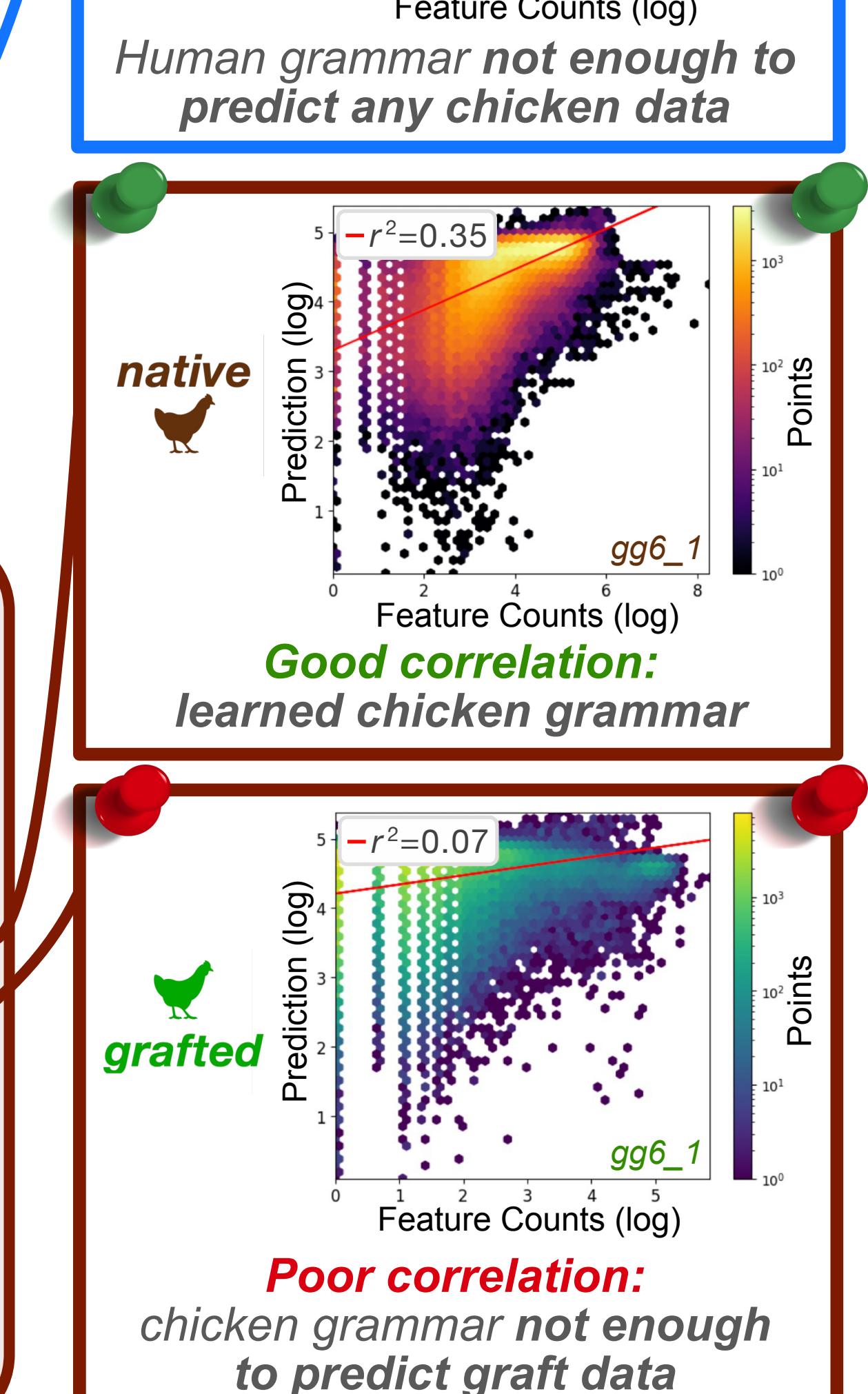
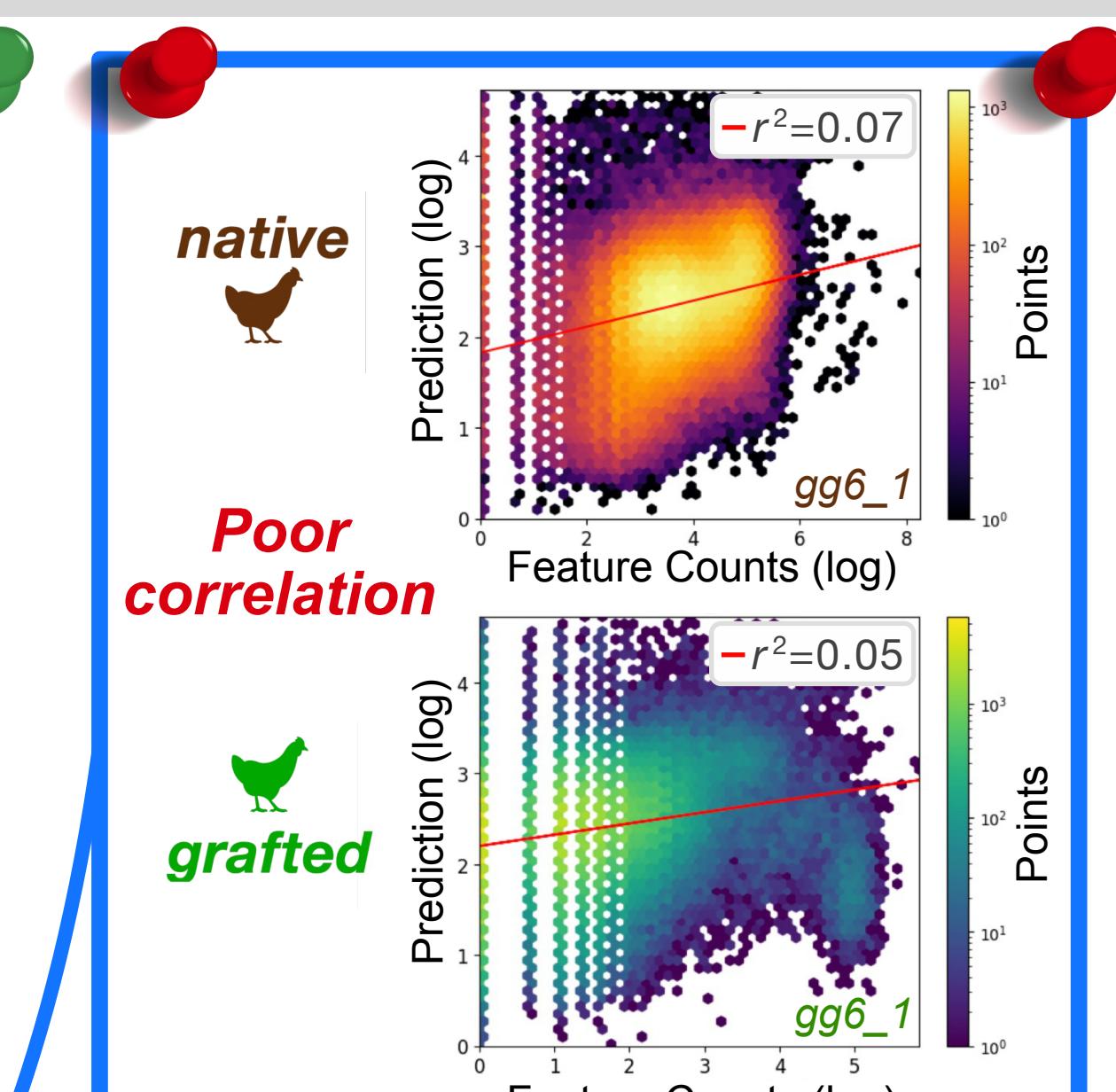
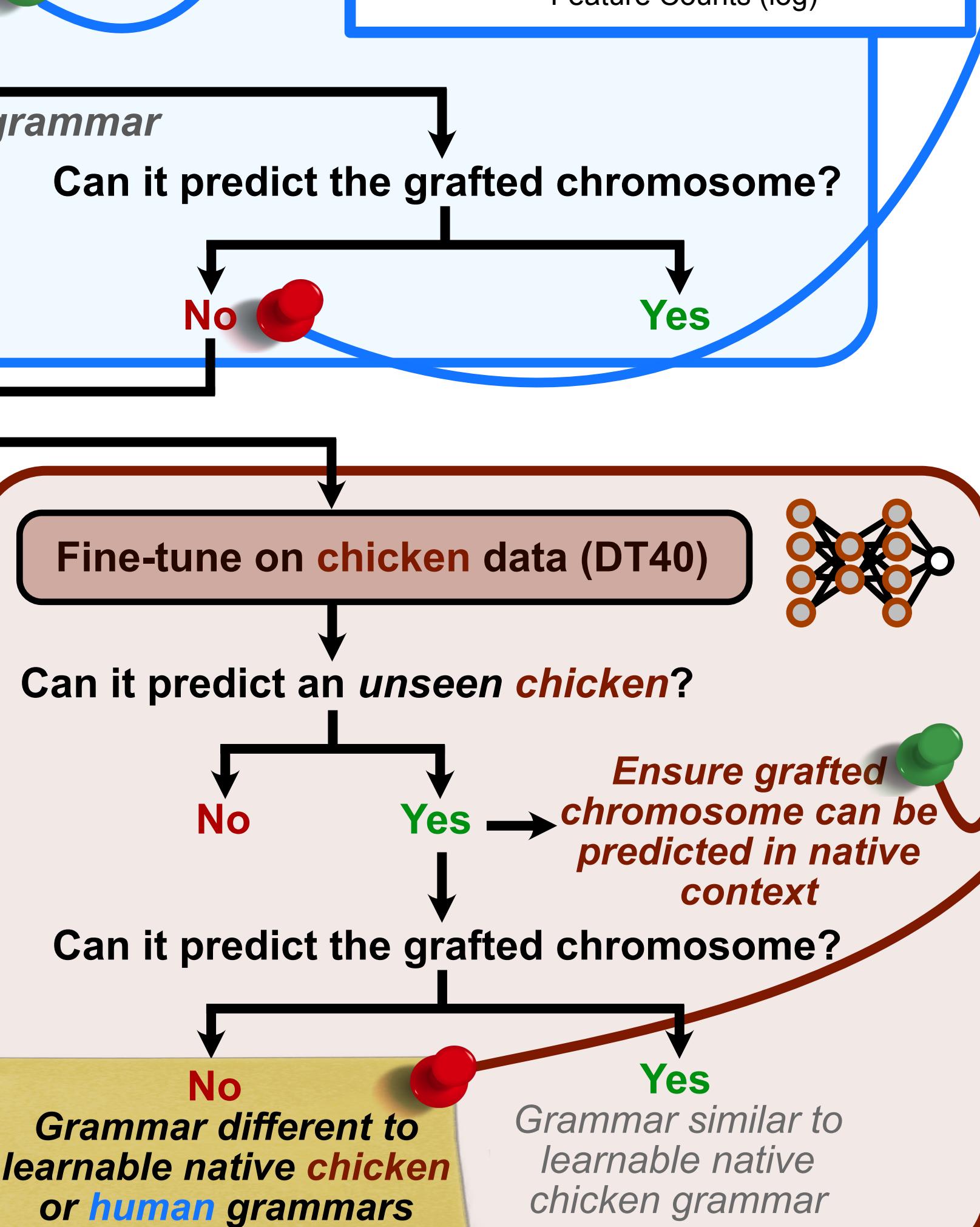
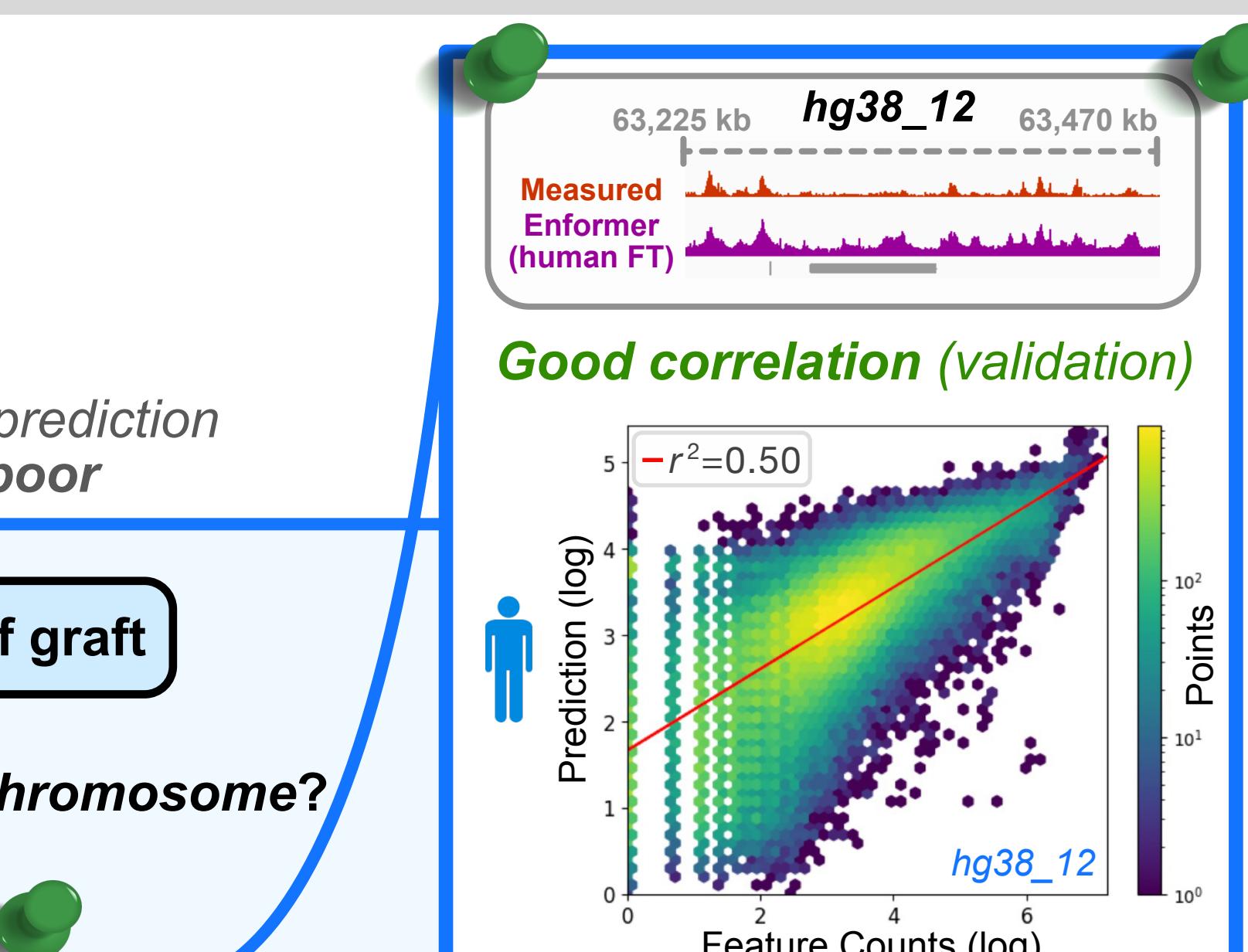
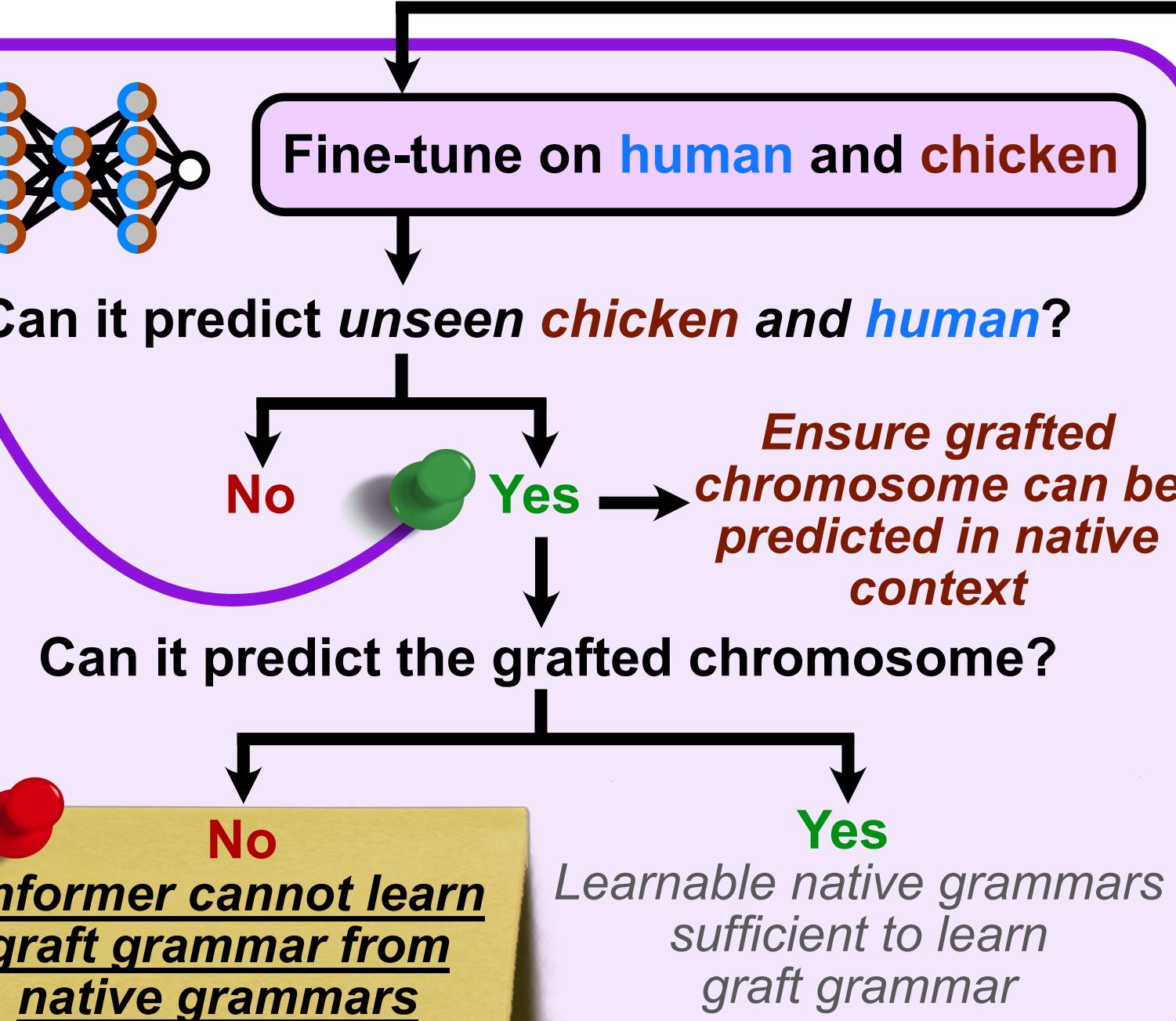


No cross-validation*, but validation loss is monitored. Trained for 500 epochs.

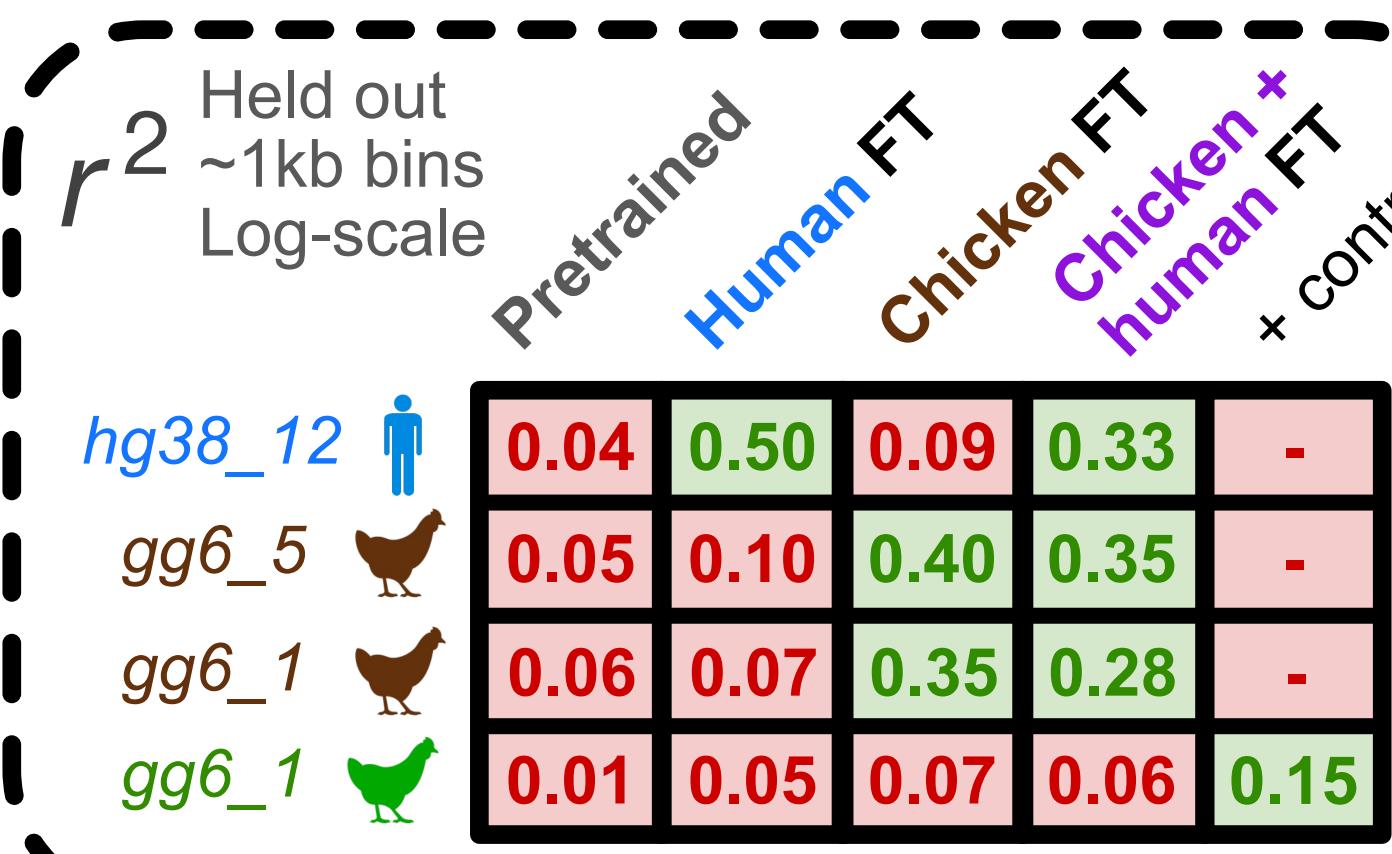
④ Enformer can learn native grammar after fine-tuning, but cannot learn grammar from graft



Poor correlation: learning both native grammars is not enough to predict graft



Poor correlation: chicken grammar not enough to predict graft data



⑤ Conclusions and future directions

- ★ After task-specific training, Enformer can learn **native regulatory grammar**
- ★ Apply workflow to other histone marks and data types

- ★ Enformer cannot learn regulatory rules of the grafted chromosome, suggesting they are neither **chicken** nor **human-like**

- ★ Integrate biological information to inspect model mispredictions

References

- Avsec, Ž., Agarwal, V., Visentin, D. et al. Effective gene expression prediction from sequence by integrating long-range interactions. *Nat Methods* **18**, 1196–1203 (2021).