# Predicting G-quadruplex Formation

**Jacob Calvert[1], Alex Kreig[1], Saurabh Sinha[2] and Sua Myong[1]**

[1]Department of Bioengineering and [2]Department of Computer Science, College of Engineering, University of Illinois at Urbana-Champaign
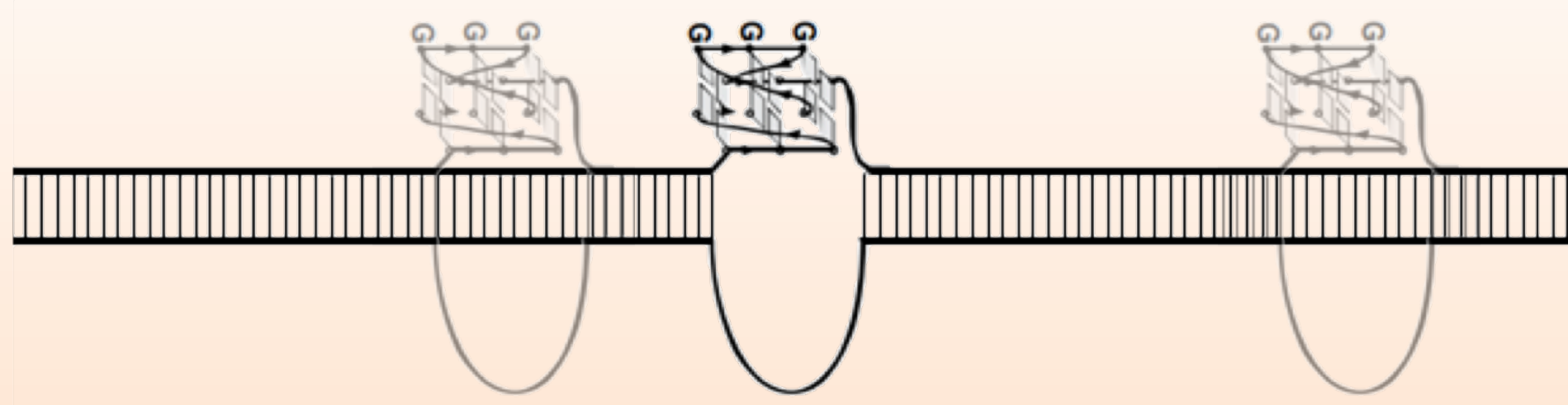
## Aims

- Model G-quadruplex folding
- Predict the folding of new sequences

## Introduction

**Q: What is a G-quadruplex (GQ)?**
**A:** *A region of guanine-rich DNA that* **can** *fold into a hill-like structure (Fig. 1a,c)*

**Q: Where are GQs found?**
**A:** *In regulatory regions like gene promoters (Fig. 1d) and telomeres (Fig. 1a,b)*

**Q: Why model GQs?**
**A:** *To better our understanding of gene regulation and motivate new disease therapies*

(TTAGGG)$_4$

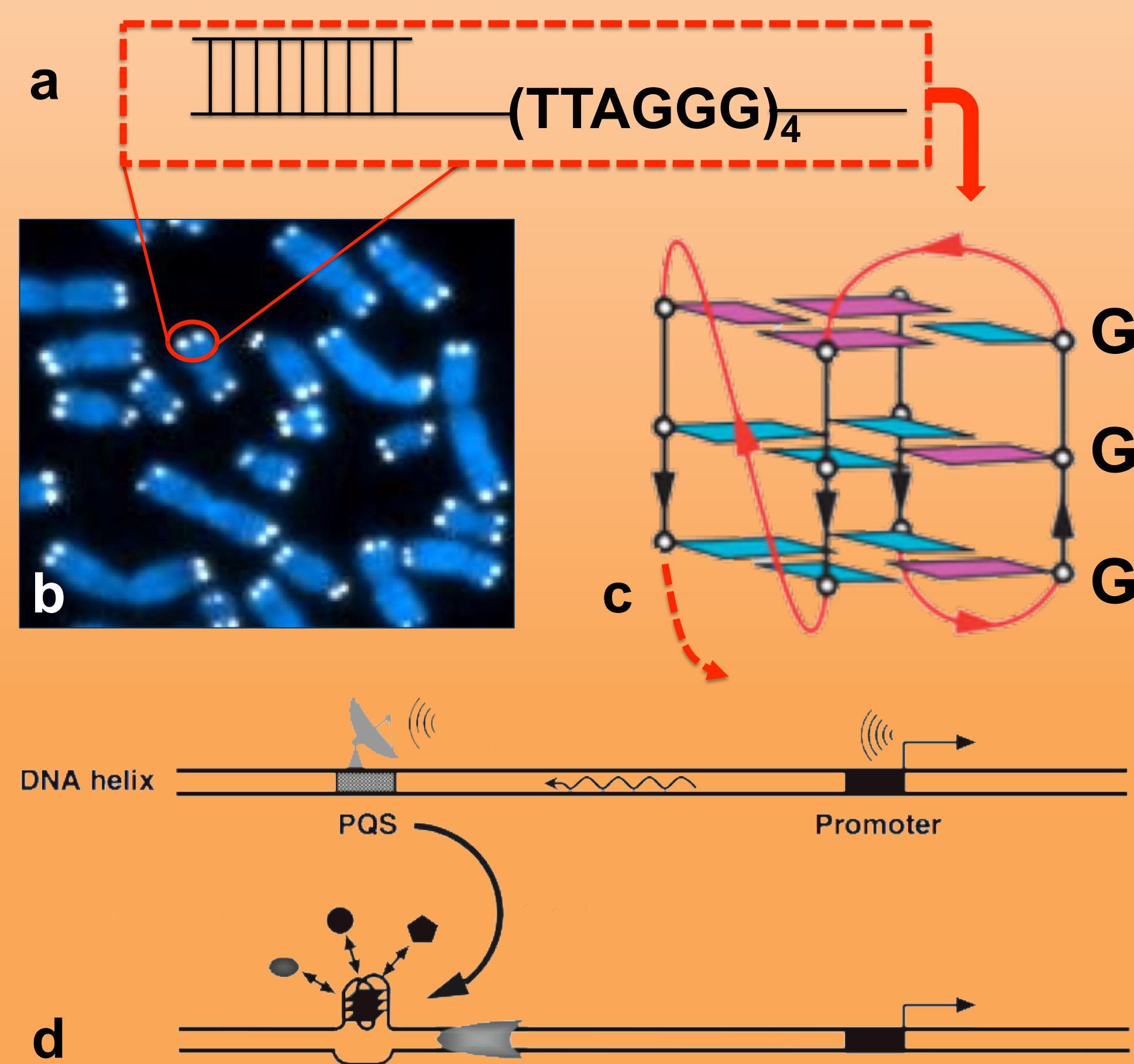G
G
G

DNA helix
PQS
Promoter

**Figure 1.** *(a) Telomere scheme. Adapted from Yildiz Lab. (b) Telomeres (white) cap the ends of chromosomes. (c) Parallel, planar interactions stabilize the GQ. Adapted from Phan AT et al. (2007). (d) Transcription-activated GQ formation in gene promoter. Adapted from Zhang C et al. (2013).*

## Data

**Q: Where did the data come from?**
**A:** *"Pull-down" experiments detected folded GQs in human cells. We compared folded GQ sequences to GQ sequence motifs that did not fold (Fig. 2a).*
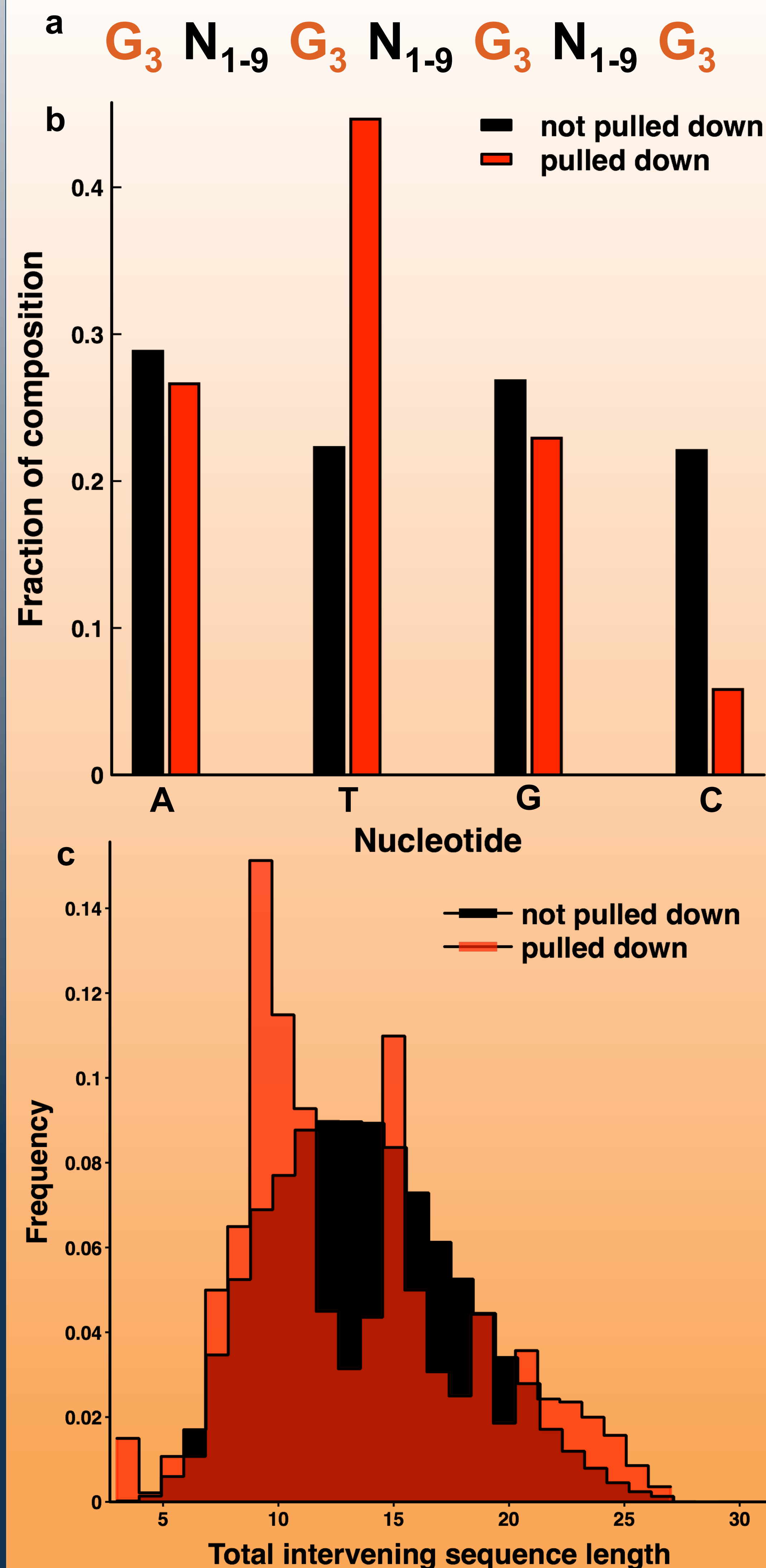
**a**  $G_3$ $N_{1-9}$ $G_3$ $N_{1-9}$ $G_3$ $N_{1-9}$ $G_3$

**b**

- not pulled down
- pulled down

Fraction of composition — Nucleotide (A, T, G, C)

**c**

- not pulled down
- pulled down

Frequency — Total intervening sequence length

**Figure 2.** *(a) GQ sequence motif.* **N** *represents any base,* **N$_{1-9}$** *are called intervening sequences. (b) Comparison of base composition and (c) the total intervening length (sum of three* **N$_{1-9}$** *loops).*

## Results

**Q: How are the pull-down data used?**
**A:** *We use a probabilistic model to detect the unique features of the pulled-down sequences. We translate the probability that a sequence folds into a score, which we call the "QPD Score."*
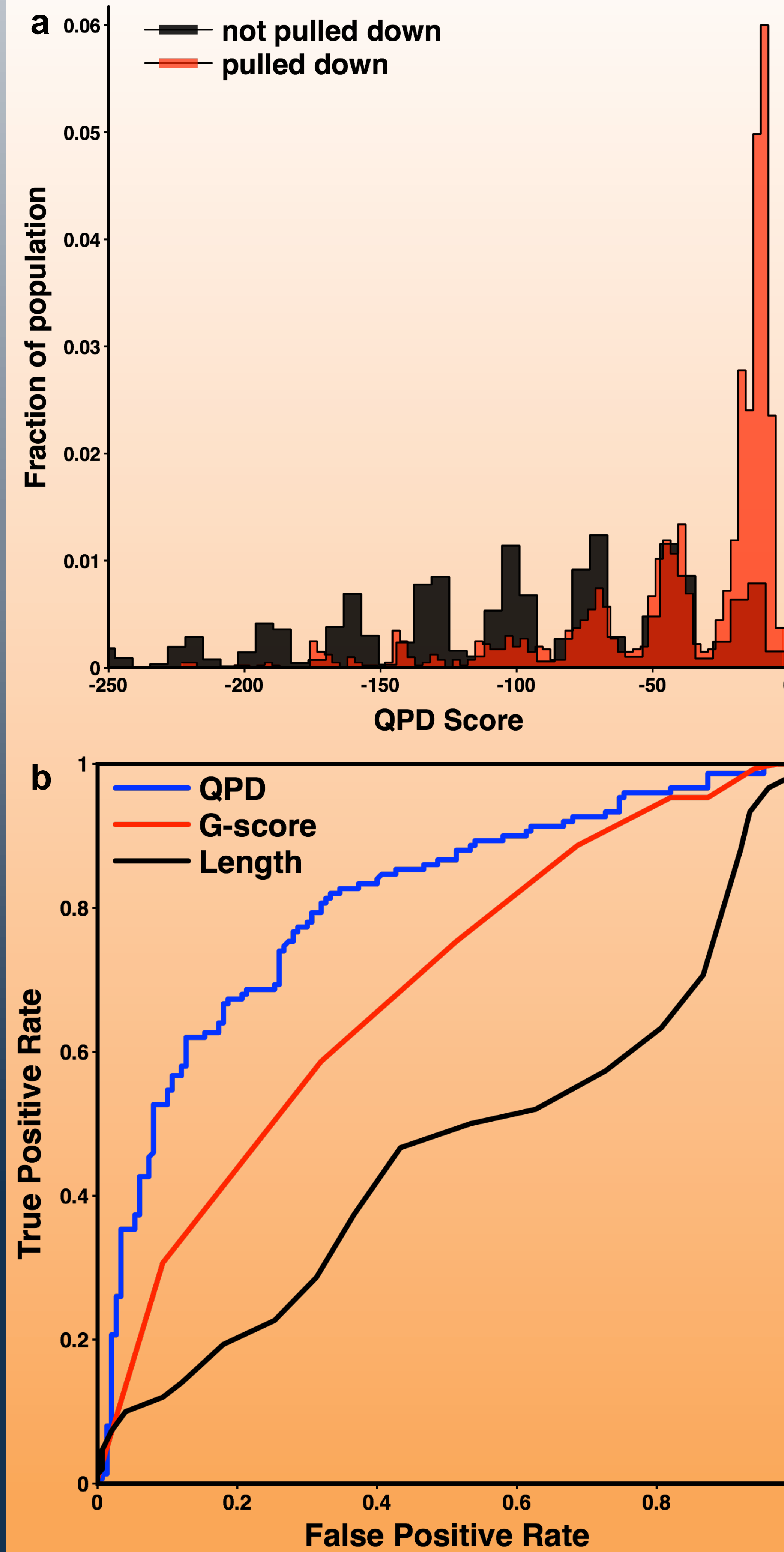
**a**

- not pulled down
- pulled down

Fraction of population — QPD Score

**b**

- QPD
- G-score
- Length

True Positive Rate — False Positive Rate

**Figure 3.** *(a) Distribution of QPD scores for all genomic GQs. (b) ROC curve comparison of QPD score against two existing methods for predicting GQ folding. Plot was generated with sequences not included in the training set.*
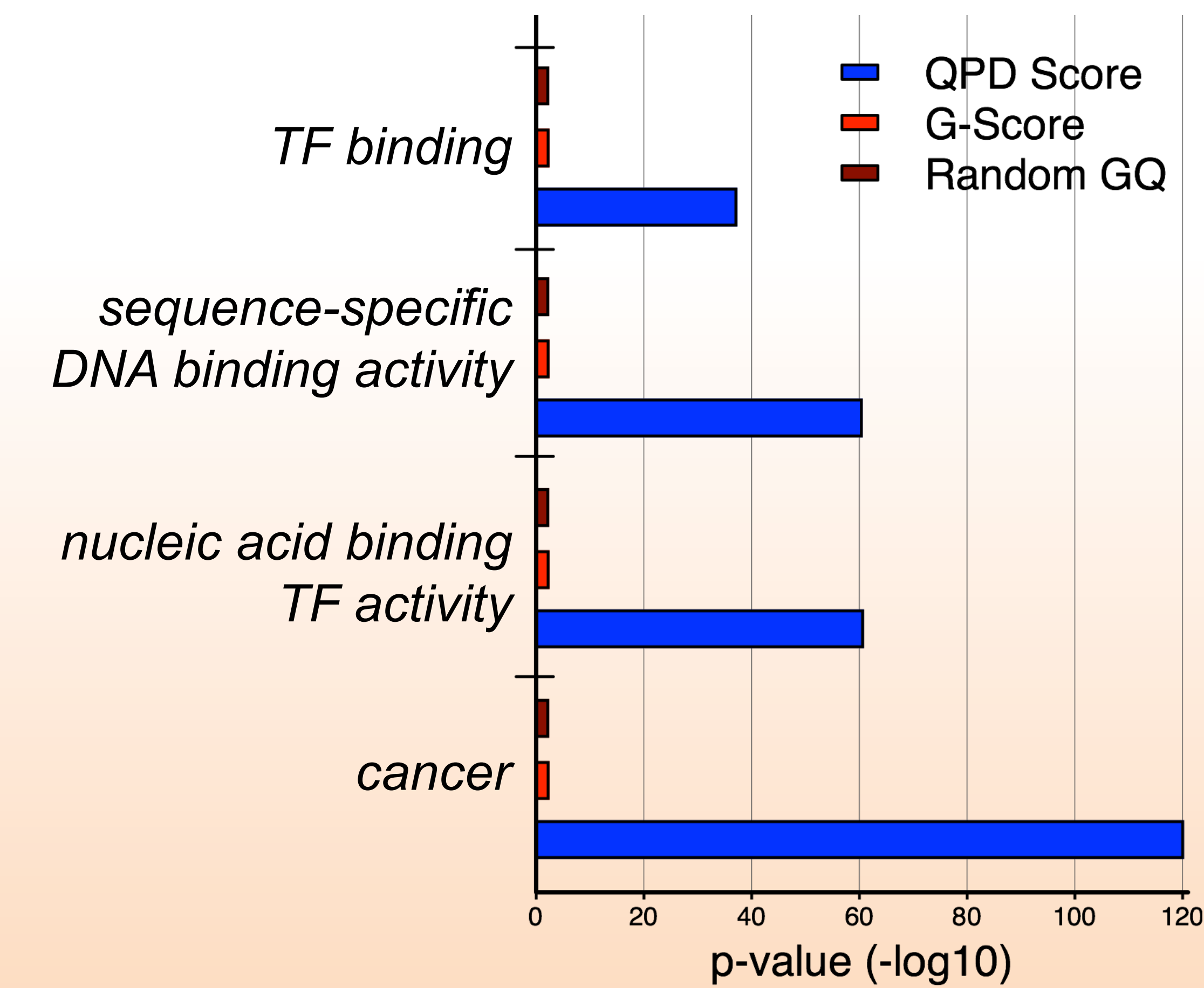
- QPD Score
- G-Score
- Random GQ

TF binding
sequence-specific DNA binding activity
nucleic acid binding TF activity
cancer

p-value (-log10)

**Figure 4.** *Selected ontologies. High QPD-scoring GQs tend to localize near genes of importance to transcription factor (TF) activity, regulation, and cancer. Sequences chosen in other ways do not.*

**Table 1.** *A selection of genes whose promoters contain GQs predicted to fold by QPD. Sequences are listed with* **G$_3$** *omitted.*

| Name | Sequence | G-Score Percentile | QPD Percentile |
|---|---|---|---|
| ABL2 | --AAGGA--A--A-- | 54 | 98 |
| RAB31 | --T--A--GTAGA-- | 69 | 99 |
| MSLN | --T--TGAA--GT-- | 69 | 99 |
| PBX1 | --AATA--GT--AGT-- | 82 | 97 |
| BCL3 & ERG | --A--A--A-- | 95 | 99 |
| WNT10A | --T--T--G-- | 95 | 99 |

## Conclusions

- Our model outperforms existing methods of GQ folding prediction.
- Highly-scoring sequences localize near genes important in regulation and disease

## Acknowledgments