

Integrative dissection of gene regulatory elements at base resolution

Zeyu Chen^{1,2,3,6}, Nauman Javed^{1,2,3,6}, Molly Moore², Jingyi Wu^{1,2,3}, Gary Sun^{1,3}, Michael Vinyard^{2,4,5}, Alejandro Collins², Luca Pinello^{2,4}, Fadi J. Najm^{2,*}, Bradley E. Bernstein^{1,2,3,7,*}

Summary

Initial submission: Received : 9/25/22

Scientific editor: Laura Zahn

First round of review: Number of reviewers: 4
Revision invited : 11/15/22
Revision received : 2/21/23

Second round of review: Number of reviewers: 4
Accepted : 3/31/23

Data freely available: Yes

Code freely available: Yes

This transparent peer review record is not systematically proofread, type-set, or edited. Special characters, formatting, and equations may fail to render properly. Standard procedural text within the editor's letters has been deleted for the sake of brevity, but all official correspondence specific to the manuscript has been preserved.

Referees' reports, first round of review

Reviewer #1:

The paper presents new experimental assays based on Cas9 to dissect the effect of gene regulatory elements. These include perturbation of entire regions, perturbation of larger binding sites, and perturbation of single base pairs. The authors focus on utilizing these assays in a single region (the CD69 locus). Using a variety of read-out subsequent to genomic manipulations, the authors identify two transcription factors in CD69 locus that have a major impact on expression changes during differentiation, and conclude that the observed opposing effect of these two TFs is explained by direct steric competition between them.

The paper has two aspects to it: describing new assays for dissecting regulatory architecture of a given genomic region, and deriving new biology through the application of these assays. I think the first component can stand on its own, however I have concerns about the conclusion the authors draw from the presented results about the steric competition between GATA3 and BHLHE40. In short, I don't think the presented results directly support the steric competition conclusions. I will summarize major and minor points below.

Major points:

- Steric interaction between GATA3 and BHLHE40 not directly supported by data. If I understand the presented data correctly, the authors do not present convincing evidence for an exclusive effect between GATA3 and BHLHE40 binding sites that are within proximity and necessary for steric competition. The results of the repression and overexpression experiments were not presented in sufficient resolution for specific locations of the two motifs. The conclusion for steric hindrance was mainly drawn from the Enformer attribution scores within the most suppressing region (sg#70). However, it is insufficient to infer and generalize from this methodology at this single region (Also see comment about Fig 3D below). Fig. 3F shows that BHLHE40 binds more across the entire region (because of the resolution) not that it binds more near the disrupted GATA3 binding location in Fig 3D (sg#70) and therefore I would suggest that steric hindrance is not the main mechanism for competition with GATA3. Figure 4E presented a statistical test to determine the effect of the distance between motifs of GATA3 and BHLHE40 to chromatin accessibility at GATA3 binding sites upon BHLHE40 overexpression. I could not find the description of the statistical test that was performed, nor any information on whether the distribution of the selected sites was accounted for in this test. Therefore, it is unclear if the significant distances are due to the sampling bias of the selected sites (i.e. more sites with close motifs in the entire set).
- Related to above, Fig 3D is confusing and hard to parse. Specifically, the predicted attribution score at BHLHE40 is low, and hard to determine the role of the other two possible motifs (TG CAG motif or GGA_G motif) with similar prediction scores. Do you have any data on the effect of base-pair changes in one of the four positions in the BHLHE40 motif to confirm that BHLHE40 is actually binding there? The CisBP and Jaspar motifs (which are consistent across experimental platforms and species) for BHLHE40 are different from the one shown. Based on the BHLHE40 motif and the weak importance scores at the interpreted binding site I wonder whether there is enough evidence to come to the conclusion that is made throughout the rest of the paper.

- When you say that you “pinpoint individual functional bases in these REs” to complement the dCas9 tiling with CBE and ABE screens, it suggests that your screen is covering most of the bases. However, if I understand correctly, most of the time only very few bases are mutated and therefore your description of the method is a bit misleading and I would suggest to clarify that this screen is not capable of producing results at base-pair resolution.
- Not clear to me that the results of Enformer are useful or necessary for the paper. Why are the Enformer results shown in Fig 1B, when they are actually ignored in selecting the regulatory regions (which is done based on ATAC signal)? One experiment to make Enformer results more convincing and relevant would be to assess the correlation between the Enformer predictions and the measured effects from the single nucleotide variations (Fig 1D).

Minor points:

1) Details of results were not sufficiently described to grasp the validity of the results:

- *“We refined these predictions using the Enformer model”*. I’d add a sentence when this is mentioned to describe what you mean by “refined”, and point to the method section describing fine-tuning.
- *“We refined these predictions using the Enformer model (Avsec et al., 2021) trained on chromatin maps and CAGE-seq data. Genomic intervals corresponding to the promoter, 3’ UTR and an RE located ~4 kb upstream of the TSS were predicted to impact CD69 transcriptional induction (Figure 1B).”*
Labels should be mentioned in the text in brackets to link the locations to figure (f.e. Promoter (RE-3, 3.1,3.2)). Introducing the RE’s in Figure 1B would be beneficial to understand the rest of the text.
- *“ we designed a library of 101 sgRNAs that tile sequences spanning RE-3 and RE-4”*
Can you describe why 101bp region is sufficient, and what would be the expected length of the perturbed region bt dCas9

2) Figure 1C:

- Did not see the legend for the colors and IDs in Fig1B at first. Might want to point to this in the figure caption.
- Do the sgRNA for RE-3 even have an effect on expression? The green line is very close to the grey control and might be within the variance of the control?
- It would be important to visualize the standard error for control and sgRNAs to illustrate that curves really differ

3) Figure 1F:

- Could you please put a brief high-level description of the method here to understand how many base changes are introduced with this system, so that the reader is able to interpret the results.
- Do these editors only edit at one position of the region that is bound by CAS9 or do they randomly edit any base in that region?
- If the edits are only at one or two positions, would you say that they are less sensitive than the dCas9 system since one base does not necessarily destroy a motif while the dCas9 system covers the entire binding site.

4) *"Notably, the CBE and dCas9 perturbations both pinpointed a ~150 bp interval within RE4 centered at sg#70 as critical for CD69 expression (Figure 2A; Chr12:9764860-9765010)."*

- Why did you not mention the second region sg#48 that you can see in 1D, E, F, 2A, and S1D here?

5) *"Several ABE hits in or near this interval also suppressed CD69 induction, but with lower fold-enrichment, potentially due to reduced effect sizes (Figure S2G)."*

- I am not sure what you mean here when you say that lower fold-enrichments were potentially caused by reduced effect sizes? How do you explain the increase for some ABE's in the region?

6) Fig 4E:

- Why would there be an effect outside close interaction range if the two factors influence each other through steric interactions?
- How did you compute the significance of the counts with a certain distance? What's the null distribution of the selected sites?
- Why did you not compare to the overall number of selected sites but instead chose the number of increased sizes for a comparison?

7) *"Sites that were repressed by BHLHE40 overexpression showed a strong enrichment for motif spacing of 0 to 3 bp, consistent with steric hindrance and competition between factors (FDR < 0.05)...."*

- In contrast to your statements, there is significant repression at 6,8,and 9, whereas there is only significant repression at position 1 and 4. Moreover, there are a significant number of negative effects at various positions and it is unclear whether the shrinking number is due to smaller effects with distance between motifs or because of the distribution of selected motifs.
- Could you please elaborate on what exactly you mean by the previously reported sites of coordinate GATA Ebox factor binding? Does BHLHE40 have a known activating effect if it is 9bp away from GATA3?

8) *Methods section:*

"Model interpretation was conducted as described at <https://github.com/deepmind/deepmind-research/blob/master/enformer/enformer-usage.ipynb>. For CAGE-seq interpretation, we calculated the gradient of the model for unstimulated Jurkat T-cells with respect to the predicted CAGE-seq signal at the CD69 promoter. This was achieved by centering a 393216 bp genomic window within the CD69 promoter (chr12:9760820-9760903)..."

- Could you briefly describe why you selected regions within the promoter and not downstream or around TSS to model expression?

"... and computing the gradient for human output head # 4831 with respect to output bins 446-450. The absolute value of the gradients were then summed in 128bp bins for coarse grain resolution (Fig 1D).

A similar approach to nominate bases contributing to RE-4 accessibility was adopted to obtain the base resolution contribution scores for the fine-tuned model corresponding to Figure S2 and 3. For this analysis, the window was centered around RE-4(chr12:9764300-9765900) and the gradient was computed with respect to output bins 442-454(Fig S1D, 3D)."

- Could you briefly describe why you used more windows than before?

9)

"The implicated interval in RE-4 is over-represented for multiple TF motifs relevant to immune function, including GATA, bHLH/Ebox, TCF, ETS and STAT (Figure 2B)."

- What does overrepresentation mean in this context? Did you do a statistical test? Could you describe how you selected (measured over-representation) these five TFs in the Methods section and point to it in the text.
 - If a statistical test was performed, it would be helpful to put a number that describes the overrepresentation in figure 2B
 - It would also be helpful to put the motifs of these TFs next to their names in Fig2B to help the reader understand what their binding sites have in common or how diverse they are.

10)

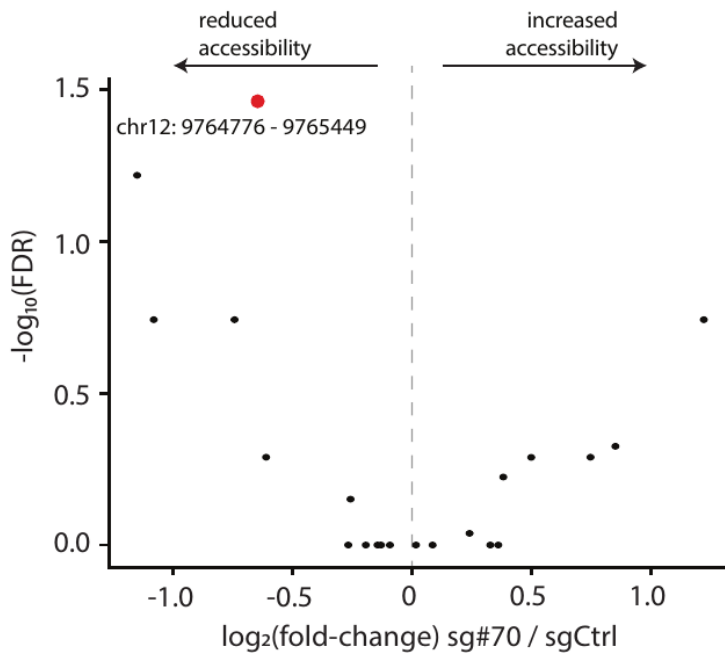
"Notably, a second top scoring interval from the CBE and dCas9 screens, centered at sg#48, showed similar TF motif enrichments (Figure 2A; Chr12:9765200-9765310)."

- You should at least mention this region earlier in the text because it shows up as early as Fig1D and it is confusing for the reader to wait for this information until now.

11)

"Fig S3D) Volcano plot depicts chromatin accessibility changes between sgCtrl and sg#70 groups within a 2mb window around RE4(red). X-axis shows log₂(fold-change) for sg #70 peaks relative to sgCtrl while Y-axis shows -log₁₀(FDR), with BH correction of p-values based on changes within 1 mb window around RE-4."

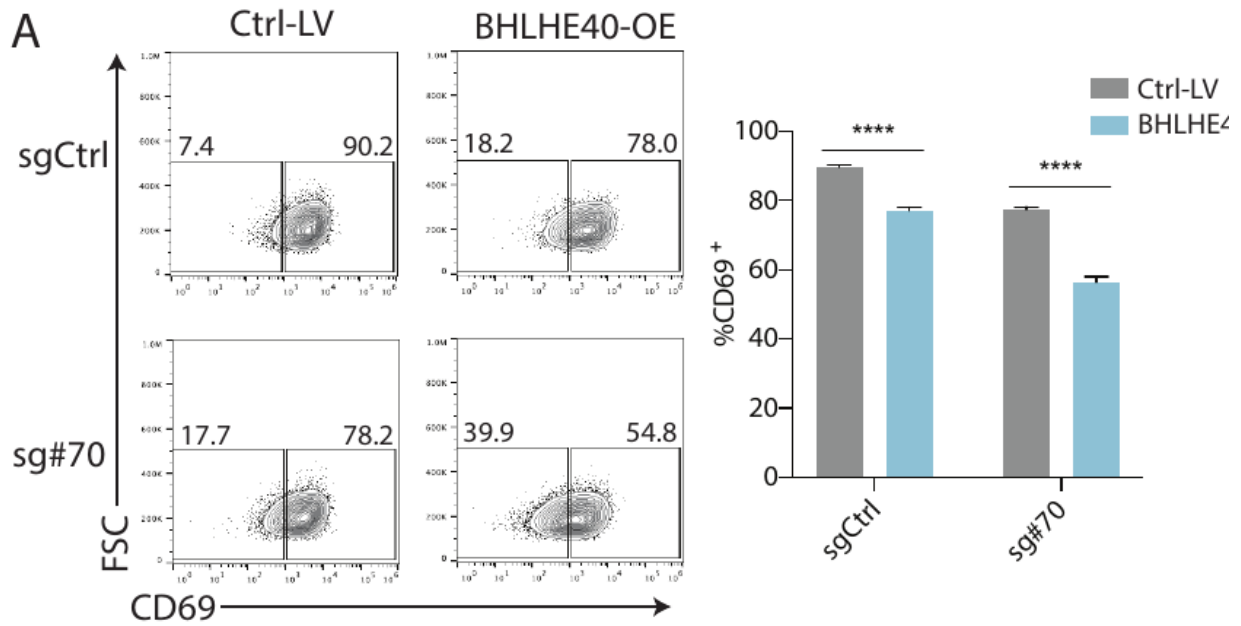
D



- Each dot represents changes over how many bases?
 - 1mb seems too large
 - What's the dot closest to the red one?
- Should be explained in more detail in Methods and in the figure caption

12)

"We found that *BHLHE40* overexpression suppressed *CD69* induction in both control and CBE-sg#70 edited Jurkat cells (Figure 4A). However, the magnitude of suppression was greater in the edited cells, potentially due to relief of *GATA* factor competition."



- Is the magnitude of suppression really greater? Both are 4-star to ctrl. How to properly compare 90.2 → 78.0 versus 78.2 → 54.8

13)

"...in our examination of the second interval identified in our dCas9 and CBE screens. Remarkably, the top base edit hit in this interval (sg#48) also incurs a C->T edit that disrupts a GATA motif flanked by a bHLH/Ebox motif (Figure 2A-2B)"

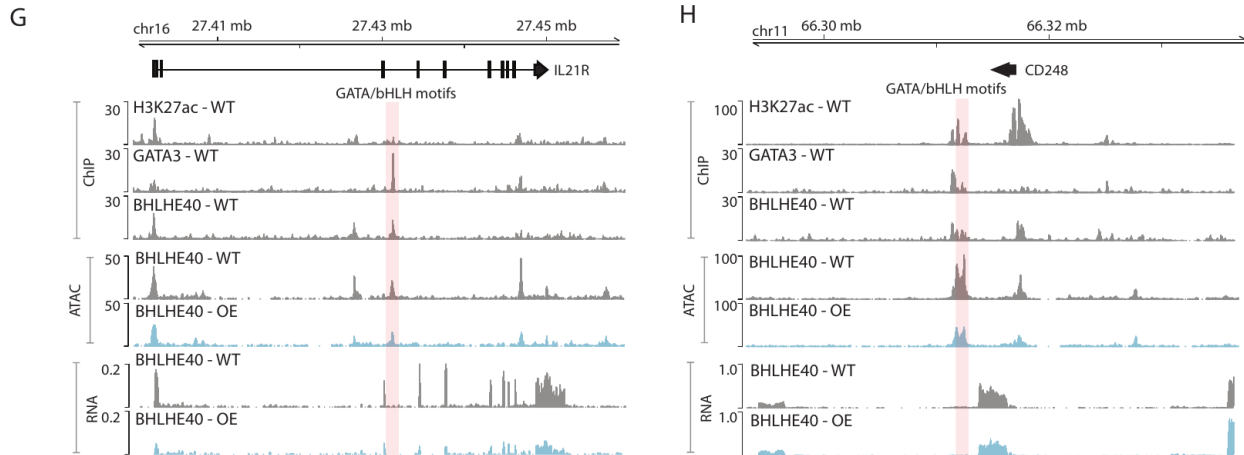
- Could you please provide an enlarged image of the region in the supplementary. Figure 2B shows a lot of GATA/BH combinations at other locations too.

14)

"We collated all GATA3 bound sites in Jurkat cells that contain a GATA motif and a bHLH/Ebox motif within the corresponding accessible site."

- Could you please explain briefly in more detail what was done here or point to the methods section where this is described.
- How did you collate these GATA3 bound sites in Jurkat cells?
- I assume you're looking at CHIP-seq from GATA or just motif scans?

15) Fig 4G,H



"These results are consistent with a general role of BHLHE40 in restraining GATA3 mediated activation at immune loci. Notably, many of the immune loci subject to opposing regulation contain elements with closely spaced GATA and bHLH/Ebox motifs (0-3 bp), consistent with a general role for competitive TF binding on T cell transcriptional programs and phenotypes (Figure 4G-4H)."

- The two presented figures do not possess sufficient resolution to confirm that the close interactions of GATA and BHLHE40 are necessary for the effect.
- Moreover, there are other positions that show way higher changes in ATAC-accessibility with overexpressed BHLHE40 compared to the regions marked in red.
- These two figures indicate that BHLHE40 is a global repressor that also represses access of GATA3 to its regulatory element, however it does not confirm that this effect takes place through steric hindrance.

Reviewer #2: Chen, et al. present a comprehensive dissection of a regulatory element that controls CD69 expression in stimulated T cells. Their approach involves an innovative combination of regulatory genomics assays, deep neural network predictions, CRISPRi, dCas9 tiling, and CRISPR-guided precision base editing. The manuscript convincingly demonstrates that a regulatory element highlighted by differential accessibility and neural network predictions (RE-4) is responsible for CD69 expression in stimulated T cells, and it further shows that GATA and E-box motif elements in RE-4 are strongly contributing to the expression response.

This manuscript represents a new and powerful combination of computational and experimental approaches, where the base-resolution regulatory predictions arising from deep neural networks are tested using base-resolution CRISPR screening techniques. I enjoyed reading the manuscript, but I remain confused and unconvinced by a couple of points.

1.a) The manuscript is a little opaque about exactly what the Enformer neural network contribution scores (i.e., Fig. 1D and Fig S1D) represent and what logic was followed when using the scores to create regulatory predictions. The Enformer scores in Fig 1D appear to be derived from a component of the model that is trained to predict unstimulated T cell CAGE data. Thus, with respect to CD69, the model should be expected to highlight regulatory elements that contribute to the *low* expression of CD69 in unstimulated T cells. It is not clear to me why these contribution scores are used to justify experimental investigation of elements that control higher CD69 expression in stimulated T cells.

1.b) According to the Methods, the Enformer model was fine-tuned (i.e., re-trained) to add components that predict ATAC-seq in resting and stimulated T cells. I believe these models produce the contribution score tracks in Fig. S1D (although this should be clarified in the figure caption). However, these models do not seem to be used for any purpose in the investigation of the RE-4 element. Are any sequence elements highlighted as differentially scoring in the stimulated vs. resting ATAC-seq Enformer models? This would appear to be a more relevant model-driven question than the unstimulated T cell CAGE model scores presented in Fig. 1D.

2) Competition between GATA3 and BHLHE40 is proposed to explain expression dynamics in stimulated T cells. This is an appealing model, but I am not fully convinced by the details. Firstly, BHLHE40 is assumed to be a universal repressor, but this is not demonstrated. There are many cases where a TF that is generally repressive can nonetheless play an activating role at a subset of sites. Indeed, BHLHE40 over-expression apparently leads to just as many activated genes as repressed genes in T cells (Fig. 4F).

Secondly, the proposed mechanism states that the observed spacing between GATA and E-box motifs is "too close to permit concurrent binding". This is not supported by any modeling or analysis. In the very least, it should be possible for a monomer of BHLHE40 to bind a half-site E-box alongside GATA3; the spacing is not that much less than the observed spacing between GATA motifs and TAL half-site E-boxes in erythroid lineages (PMID: 26503782).

Thirdly, the sg#70 CBE edits result in expanded BHLHE40 ChIP-seq signal at RE-4, reduced H3K27ac at RE-4, and reduced expression of CD69 in the context of BHLHE40 over-expression. However, the sg#70 CBE should equally affect both the GATA motif and the neighboring E-box (Fig. 3D). Thus, the sg#70 results cannot by definition support the presented model.

Overall, it is not clear why the competition model was favored over activating cooperation between GATA3 and a bHLH TF. It would be interesting to see what effect a BHLHE40 knock-down would have on CD69 expression.

Reviewer #3: This study by Chen et al. described a method of integrative analysis combining epigenetic perturbations, base editing, and deep learning models for dissecting gene regulatory elements at base resolutions. Taking CD69 gene locus as an example, this analysis identified an artificial C-to-T variant that suppress CD69 expression. This C-to-T base edit was shown to ablate a GATA3 binding site and eliminate GATA3 binding at this site in Jurkat cells. The authors further suggested a potential mechanism of binding competition between GATA3 and BHLHE40 in regulating inducible immune genes and T cell state. While the topic of determining the functions and sequence determinants of cell type-specific regulatory elements is of great interest, this study as presented could not demonstrate that this integrative analysis and machine learning proposed by the authors can help answer this big question.

1. Why does the authors choose to use CRISPRi to test the functional impact of the promoter (RE-3), the putative upstream RE (RE-4) and two other sites (Fig. 1B)? Is this decision purely based on the differential ATAC peak analysis? Did Enformer prediction help here? Even though the authors claimed that "the Enformer model predicted that a specific ~170 bp sequence interval within RE-4 is most critical for CD69 regulation", without this prediction, dCas9 tiling and CBE tiling can also give the same conclusion. The machine learning here did not help narrow down candidates or reduce the amount of screen that need to be performed.

After CBE tiling and subsequent sequencing pinpoint C-948 to be critical in regulating CD69 expression, a transcription factor motif analysis can identify its location in a GATA factors binding site. How does the Enformer prediction help in identifying the responsible factor?

The machine learning or Enformer prediction does not seem to give new information other than those already provided by other analysis, nor seem to save effort in screening.

2. Why is that the authors analyzed BHLHE40, but not BHLHE22 or ARNT2 for their potential role in regulating CD69 expression, while all of them are strongly induced in stimulated Jurkat cells? Even though CBE-sg#70 eliminate GATA3 binding at RE-4, the gain of "broader BHLHE40 binding over RE-4" is not convincing. Could BHLHE22 or ARNT2 show a stronger gain of binding?

3. To support the conclusion that GATA3 and BHLHE40 compete at RE-4 to regulate CD69 expression, could the authors show the effect of BHLHE40 depletion on GATA3 binding at RE-4 and vice versa? How does GATA3-OE and sgBHLHE40 influence CD69 expression?

4. Could the authors explain what the black dots in Fig. S3D indicates, please? Besides the red dot, a few black ones also show significant change in accessibility. Why does the authors state that "we did not observe any other accessibility changes in the CD69 locus or neighboring genomic regions"?

Reviewer #4: In the manuscript "Integrative dissection of gene regulatory elements at base resolution"

Chen et al. present a framework to identify and validate the impact of single base edits on regulatory element function and target gene expression. For their study they focus on the CD69 gene in T-cells which is induced in T-cell activation. Based on ATAC-seq and machine learning they identify a ~170bp interval within an enhancer predicted to regulate CD69 gene expression. Using base editing they highlight a C-to-T transition that leads to loss of GATA3 binding accompanied by loss of chromatin accessibility. They report that this leads to increased binding of the repressor BHLHE40. They also indicate that this GATA3/BHLHE40 competition is a general genome-wide process during immune cell response and T-cell activation.

The current study presents a nice, generalizable experimental and computation framework to characterize the regulatory landscape of an individual locus and insight into potential competition between activators and repressors. The paper is well written and easy to follow with clear experimental and conceptual rationale and execution.

Overall, the study provides potentially very interesting insight, but in its current form several aspects of the study seem preliminary.

I have several points that would need to be addressed:

1) The advantage of the integration of ATAC and Enformer compared to using the summit or TF footprinting in differential ATAC-seq peaks is not entirely clear to me? Based on both the Enformer and differential accessibility data presented in 1B, would one not predict RE3 to have the strongest effect? What about sites that have strong Enformer signal, e.g. the two peaks in Enformer signal track next to RE2?

2) Since the promoter is required for gene expression and indeed are often used as positive controls in screens it is surprising to see only mild effect when targeting the promoter proximal RE3. In addition to focusing on CD69 protein levels, ATAC-seq data (or H3K27ac) would be helpful to show that chromatin accessibility and activity of the respective REs is lower after infection with guide RNAs. This would provide a more direct readout and enable a mechanistic link that indeed RE activity is linked to expression.

3) Could the modest effect of sg#70 and overexpression of BHLH40 on chromatin accessibility, e.g. Fig. 4B and C might be because silencers are also accessible? H3K27ac or H3K27me3 levels might be better chromatin marks to study the competition of activator and repressor. For n=2 experiments please display individual values. The statistics using uncorrected p-value despite having whole genome data with thousands of peaks seems strange, particularly since there are many more pronounced sites as shown in panel D of the same figure.

4) Sg#70 seems to have the most drastic effect on recruitment of BHLHE40 to the promoter or proximal RE of CD69, whereas the signal at RE4 is comparable and BHLHE40 is already bound in the WT (Fig.3F). If GATA3 as stated prevents binding of BHLH40 why is there no clear difference in binding at RE4? How is the binding pattern of BHLHE40 upon infection with Sg#70 and overexpression of BHLHE40?

5) BHLHE40 is described as transcriptional repressor, but Fig 4F illustrates up and downregulated genes linked to BHLHE40 REs. How do the authors explain this? Can BHLHE40 context specific activate or repress? How were target genes defined? Is this a mix of direct and indirect effects? Could these be resolved?

6) The authors propose that GATA3 and BHLHE40 competition impacts transcriptional responses globally. I think to show this it would be important not only to show changes after overexpression of BHLHE40 since this can lead to artifacts. It would be better to integrate RE activity, binding of GATA3 and BHLHE40 by ChIP-seq and RNA-seq (several of these datasets are already generated as part of this study) in control and stimulated T-cells. This could directly

address several questions of competition, for example for genes that are higher expressed in activated T-cells do they gain GATA3, loose BHLHE40 at REs or does the relative ratio of binding change? How many sites are only bound by GATA3 or BHLHE40 and how many are co-bound? What about the associated H3K27ac and accessibility levels? And similar analysis could be done for downregulated genes. How many genes are dependent on the competition in this model?

Authors' response to the first round of review

Response to reviewer comments

Summary of major additions and revisions:

We thank the reviewers for their enthusiasm and comments on our manuscript. Their feedback prompted additional experimentation and new analyses, which greatly improved our study. We highlight below key points of impact and major additions in the revision. This is followed by point-by-point

responses to all reviewer comments.

Key points of impact and major additions for the revision:

1. Incorporation of deep learning model predictions. We expanded our use and evaluation of the Enformer model. First, using new CRISPRi data generated for the revision, we evaluated the predictive power at the level of regulatory elements. These new data reveal strong concordance between the predictions and the experimental data (revised Figure 1B, 1C, S1C, S1D). Second, we systematically assessed the base level correlation between the Enformer gradient signal and base editor tiling results. We find reasonable concordance between the respective approaches (revised Figure S2H-S2I). We also expand our use of Enformer, fine tuning the model for differential accessibility between resting and stimulated Jurkat cells and identifying motifs and TF interactions that help explain top scoring guides.

We recognize that many of the biological findings in our paper could have been achieved through combination of the epigenomic data and the systematic CRISPRi, dCas9 and base editing perturbations, without the Enformer component. However, given the increasing power and uptake of deep learning tools, we believe that the incorporation and evaluation of these Enformer data are an important contribution of our study.

2. Revised model and softened claims on TF interactions. In response to reviewer comments and new data collected for the revision, we revised our model regarding the TF interactions that underlie the potent sg#70 base edit. First, we identified additional partial e-box motifs in the region, including one site located 9 bp from the GATA motif (revised Fig. 3D/S3H). It completes a potential GATA:TAL1 binding site. TAL1 is a bHLH factor and known GATA3 partner. We confirmed by ChIPseq

that TAL1 binds the site along with GATA3 in Jurkat cells (revised Fig. 3F). Importantly, mutation of the GATA motif (sg#70 C->T edit) ablates binding of both TFs, consistent with cooperative TAL1-GATA3 binding to RE-4 in wild-type Jurkat cells.

ChIP-seq also revealed a diffuse increase in BHLHE40 binding across a region encompassing RE-4 and the CD69 promoter in the base-edited Jurkat cells. As noted by the reviewers, the diffuse binding suggests the importance of other features, beyond the motif highlighted in our original submission. Indeed, several additional (partial) e-box motifs in the region could contribute to BHLHE40 recruitment. A repressive role for BHLHE40 is supported by additional data showing that BHLHE40 over-expression decreases RE-4 accessibility and acetylation, and increases H3K27 methylation (revised Fig. 4C-4D), and lowers CD69 expression (Fig. 4A).

However, we recognize that our data do not prove that the opposing functions of GATA3/TAL1 and BHLHE40 are mediated by steric hindrance between the factors.

With further genetic validations (revised Figure 4B, S4B-S4C), we therefore simplified our model to state that the base edit displaces the GATA3/TAL1 complex, reducing the activity of RE-4 and allowing BHLHE40 and potentially other repressive factors to bind across the locus. We refer to

the possibility of steric hindrance, but add that this would not be sufficient to explain the relatively widespread changes in factor binding (revised section “Global interplay between BHLHE40 and GATA3 in T-cell responses” of revised text; Revised Fig. 4E-4F, S5).

3. We have addressed all remaining reviewer comments with new data, analyses and clarifications through the text. We thank the reviewers for their careful attention to our study.

Response to Reviewers

Reviewer #1: The paper presents new experimental assays based on Cas9 to dissect the effect of gene regulatory elements. These include perturbation of entire regions, perturbation of larger binding sites, and perturbation of single base pairs. The authors focus on utilizing these assays in a single region (the CD69 locus). Using a variety of read-out subsequent to genomic manipulations, the authors identify two transcription factors in CD69 locus that have a major impact on expression changes during differentiation, and conclude that the observed opposing effect of these two TFs is explained by direct steric competition between them.

The paper has two aspects to it: describing new assays for dissecting regulatory architecture of a given genomic region, and deriving new biology through the application of these assays. I think the first component can stand on its own, however I have concerns about the conclusion the authors draw from the presented results about the steric competition between GATA3 and BHLHE40. In short, I don't think the presented results directly support the steric competition conclusions. I will summarize major and minor points below.

We appreciate the positive comments about our approach, and for constructive feedback regarding both our computational analyses and hypothesized model of GATA3/BHLHE40 interaction. In brief, we have revised and softened our model, downplaying claims of steric opposition in favor of more general TF antagonism, in accordance with review comments. Please see detailed responses below.

Major points:

1) Steric interaction between GATA3 and BHLHE40 not directly supported by data. If I understand the presented data correctly, the authors do not present convincing evidence for an exclusive effect between GATA3 and BHLHE40 binding sites that are within proximity and necessary for steric competition. The results of the repression and overexpression experiments were not presented in sufficient resolution for specific locations of the two motifs. The conclusion for steric hindrance was mainly drawn from the Enformer attribution scores within the most suppressing region (sg#70). However, it is insufficient to infer and generalize from this methodology at this single region (Also see comment about Fig 3D below). Fig. 3F shows that BHLHE40 binds more across the entire region (because of the resolution) not that it binds more near the disrupted GATA3 binding location in Fig 3D (sg#70) and therefore I would suggest that steric hindrance is not the main mechanism for competition with GATA3.

We now have extensive computational and experimental analyses for the revision. The new data provide several key insights into TF motifs and binding patterns in the vicinity of the sg70 edited base, allowing us to refine our mechanistic models.

First, there are several additional partial e-box motifs in the region, including one site located 9 bp from the GATA motif (revised Fig 3D). Our fine-tuned Enformer model highlighted in particular this partial e-box (revised Fig 3D, S3G, methods section “Enformer predictions and finetuning”).

In combination with the GATA motif, it completes a potential GATA:TAL1 binding site.

TAL1 is a bHLH factor and known GATA3 partner implicated in T-cell biology¹. We confirmed by ChIP-seq that TAL1 binds the site along with GATA3 in Jurkat cells (Fig. 3F). When we next tested binding in edited Jurkat cells, we found that mutation of the GATA motif ablates binding of both TFs, consistent with cooperative TAL1-GATA3 binding to RE-4 in wild-type Jurkat cells.

ChIP-seq also revealed a diffuse increase in BHLHE40 binding across RE-4 and the 5' portion of CD69 in the base-edited Jurkat cells. As noted by the reviewer, the diffuse binding suggests the importance of other features, beyond the motif highlighted in our original submission.

Indeed, there are several additional (partial) e-box motifs in the region that could contribute to BHLHE40 recruitment. The functional importance of BHLHE40 over the site is further supported by our data showing that BHLHE40 over-expression decreases RE-4 accessibility and acetylation, and increases RE-4 H3K27 trimethylation (Fig. 4C-4D), and lowers CD69 expression in Jurkat cells (Fig. 4A).

We have revised our interpretation of the sg#70 base edit. We simplify our model to state that the base edit displaces the GATA3/TAL1 complex, reducing the activity of RE-4 and allowing BHLHE40 and potentially other repressive factors to bind across the locus. We refer to the possibility of steric hindrance, but add that this would not be sufficient to explain the relatively widespread changes in factor binding (revised results section “Global interplay between BHLHE40 and GATA3 in T cell responses”; Fig. 4E-4F). We hope that the revised and softened model addresses these points raised by the reviewer.

2) Figure 4E presented a statistical test to determine the effect of the distance between motifs of GATA3 and BHLHE40 to chromatin accessibility at GATA3 binding sites upon BHLHE40 overexpression. I could not find the description of the statistical test that was performed, nor any information on whether the distribution of the selected sites was accounted for in this test. Therefore, it is unclear if the significant distances are due to the sampling bias of the selected sites (i.e. more sites with close motifs in the entire set).

This figure panel on motif distance has been removed in the revision. We have confirmed that the remaining figure panels include appropriate statistics.

3) Related to above, Fig 3D is confusing and hard to parse. Specifically, the predicted attribution score at BHLHE40 is low, and hard to determine the role of the other two possible motifs (TGCAG motif or GGA_G motif) with similar prediction scores.

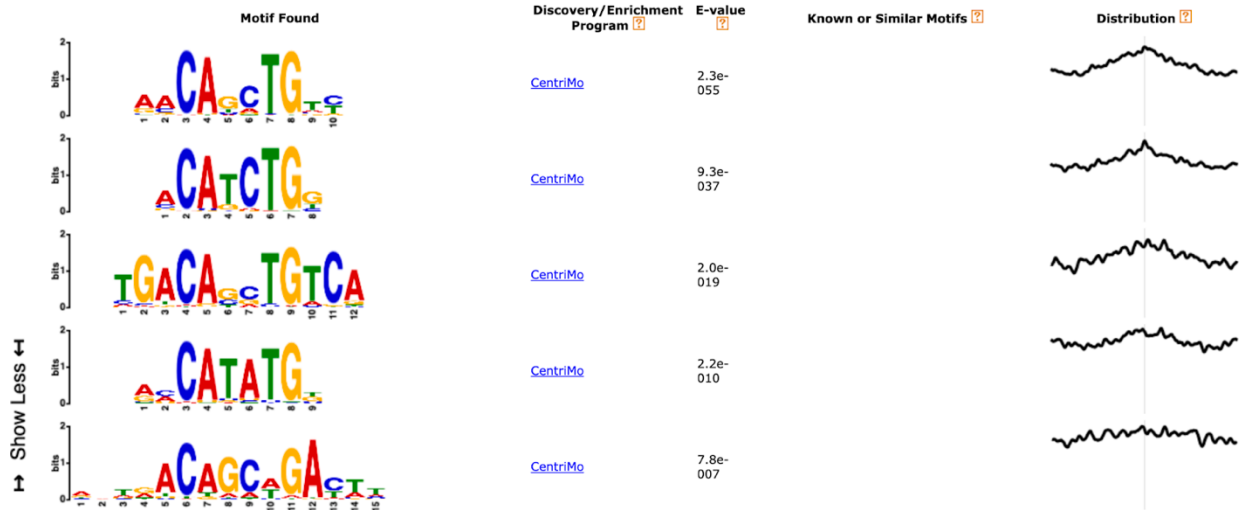
We have revised Fig. 3D for clarity. The figure now shows base contributions from the updated finetuned model (revised methods), and we have adjusted the scale and annotated these relevant motifs in the vicinity of the sg#70 edit. These include the GATA motif and the partial e-box motif that the reviewer points out, which together form a GATA:TAL1 binding site. The central e-box with the previously proposed role in BHLHE40 recruitment is also evident (though we clarify in the text that multiple e-box sites in the regions likely contribute to BHLHE40 recruitment). We do not know the significance of the GGA sequence that the reviewer points out. We caveat in the revised discussion that we do not exhaustively identify all motifs and binding sites relevant to RE-4 regulation.

4) Do you have any data on the effect of base-pair changes in one of the four positions in the BHLHE40 motif to confirm that BHLHE40 is actually binding there?

Given the reviewer’s prior comments on the diffuse binding of BHLHE40 and our identification of several additional candidate binding motifs, we have reduced emphasis on this single BHLHE40 motif (see above responses).

5) The CisBP and Jaspar motifs (which are consistent across experimental platforms and species) for BHLHE40 are different from the one shown.

We highlight the generic E-box motif (CANNTG) instead of the canonical JASPAR motif because bHLH factors also associate with a different class of E-box motifs. Consistently, our BHLHE40 ChIP (revised Fig. 3F, 4D) show central enrichment for both the expected CACGTG and the CANNTG ebox motif.



Reviewer Figure 1.1: Generic e-box (CANNTG) variants centrally enriched within BHLHE40 peaks called from Ctrl-BHLHE40 ChIP data using the XSTREME tool, plotted in Revised Figure 4D.

6) Based on the BHLHE40 motif and the weak importance scores at the interpreted binding site I wonder whether there is enough evidence to come to the conclusion that is made throughout the rest of the paper.

See above clarification of our revised model.

7) When you say that you “pinpoint individual functional bases in these REs” to complement the dCas9 tiling with CBE and ABE screens, it suggests that your screen is covering most of the bases. However, if I understand correctly, most of the time only very few bases are mutated and therefore your description of the method is a bit misleading and I would suggest to clarify that this screen is not capable of producing results at base-pair resolution.

The reviewer is correct that the fraction of bases that we can target is limited by the existence of nearby NGG PAM sites. We now clarify in the revised discussion that we cover ~28% of cytosines or guanines in RE-3 and RE-4 using CBE.

8) Not clear to me that the results of Enformer are useful or necessary for the paper. Why are the Enformer results shown in Fig 1B, when they are actually ignored in selecting the regulatory regions (which is done based on ATAC signal)? One experiment to make Enformer results more convincing and relevant would be to assess the correlation between the Enformer predictions and the measured effects from the single nucleotide variations (Fig 1D).

The reviewer makes a great point about evaluating the Enformer predictions. We address this point with two additions to the revision:

1. The revision now includes additional CRISPRi experiments that target the other regions in the CD69 locus that were nominated by Enformer (revised Fig. 1B). The additional targets include one intronic region (guides i.1,i.2), the CD69 promoter (prom.1, prom.2), and an additional distal element, RE-5 (5.1,5.2). We exclude a fourth Enformer-nominated site as it overlaps a CD69 exon. The CRISPRi guides targeting RE-4 and the CD69 promoter have the strongest impact on CD69 expression (revised Fig. 1C, S1C-S1D). More broadly, the results reveal that the measured effect of each guide correlates well with the Enformer gradient prediction for the corresponding 2kb window (Fig. S1D).

2. As suggested by the reviewer, we also assessed the correlation between the Enformer predictions and the measured effect of artificial BE-induced variants on CD69 expression. We find that base positions that score in Enformer have a greater likelihood of impacting CD69 in the base editing experiment (revised Fig. S2H-S2I).

We thank the reviewer for prompting these additional experiments and analyses. We do think these data are important and relevant here, given the critical goal of developing computational models able

to predict the function of ALL regulatory bases. With this goal in mind, the additional data provide a more robust evaluation of the strengths and limitations of the deep learning model, and thus increase the impact of our study.

Minor points:

1) Details of results were not sufficiently described to grasp the validity of the results:

“We refined these predictions using the Enformer model”. I’d add a sentence when this is mentioned to describe what you mean by “refined”, and point to the method section describing fine-tuning.

We have added a clearer explanation of how we fine-tuned Enformer in the main text (revised results section “Resolving functional bases within immune regulatory element”, paragraph 2) and expanded the methods section (revised methods section “Enformer predictions and finetuning”, paragraph 1).

“We refined these predictions using the Enformer model (Avsec et al., 2021) trained on chromatin maps and CAGE-seq data. Genomic intervals corresponding to the promoter, 3’ UTR and an RE located ~4 kb upstream of the TSS were predicted to impact CD69 transcriptional induction (Figure 1B).” Labels should be mentioned in the text in brackets to link the locations to figure (f.e. Promoter (RE-3, 3.1,3.2)). Introducing the RE’s in Figure 1B would be beneficial to understand the rest of the text.

We have reorganized Figure 1B and 1C and clarified the labels as recommended. The interrogated elements are much clearer now. We thank the review for the suggestion.

“we designed a library of 101 sgRNAs that tile sequences spanning RE-3 and RE-4”

Can you describe why 101bp region is sufficient, and what would be the expected length of the perturbed region bt dCas9

We apologize for any lack of clarity. The library of 101 sgRNAs covers ~2 kb of sequence across RE-3 and RE-4 since each guide is spaced by ~20 bp. The revised discussion more clearly caveats that we do not achieve full coverage of these regions due to limitations associated with the NGG PAM requirement.

2) Figure 1C: Did not see the legend for the colors and IDs in Fig1B at first. Might want to point to this in the figure caption. Do the sgRNA for RE-3 even have an effect on expression? The green line is very close to the grey control and might be within the variance of the control? It would be important to visualize the standard error for control and sgRNAs to illustrate that curves really differ. We have added a more prominent legend for the colors in revised Fig. 1B and indicated this in the figure caption. The sgRNAs for RE-3 do have an effect on expression as demonstrated in revised Fig. 1C. We have added data from each triplicate to the revised Fig. 1C.

3) Figure 1F:

Could you please put a brief high-level description of the method here to understand how many base changes are introduced with this system, so that the reader is able to interpret the results. Do these editors only edit at one position of the region that is bound by CAS9 or do they randomly edit any base in that region? If the edits are only at one or two positions, would you say that they are less sensitive than the dCas9 system since one base does not necessarily destroy a motif while the dCas9 system covers the entire binding site.

We have added a brief description of the system to the legend of Fig. 1F. The CBE edits any cytosine (on either strand) 2-11 base pairs opposite the NGG PAM sequence. We observe some preference for positions in a central 2 to 8 base window, supported by another study using the same CBE construct. We estimate that our library of 101 sgRNAs converts ~28% of cytosines within RE-3/RE-4 to thymines. The reviewer makes a good point that editors are likely less sensitive than dCas9 in cases where a motif is not edited/destroyed. This was part of our rationale to use both systems.

4) “Notably, the CBE and dCas9 perturbations both pinpointed a ~150 bp interval within RE4 centered at sg#70 as critical for CD69 expression (Figure 2A; Chr12:9764860-9765010).”

Why did you not mention the second region sg#48 that you can see in 1D, E, F, 2A, and S1D here?

We now include additional data showing a head-to-head comparison of sgCtrl, sg#70, and sg#48 in revised Fig. S3C. Although sg#48 clearly suppresses CD69 suppression (as the reviewer notes from other figures), it has a weaker phenotype than sg#70.

5) "Several ABE hits in or near this interval also suppressed CD69 induction, but with lower foldenrichment,

potentially due to reduced effect sizes (Figure S2G)." I am not sure what you mean here when you say that lower fold-enrichments were potentially caused by reduced effect sizes? How do you explain the increase for some ABE's in the region?

We observed lower signal-to-noise in our ABE data compared to CBE data. Further study is needed to evaluate whether this reflects lower editing efficiency or decreased likelihood for A to G edits to affect regulatory element function. To improve sentence clarity, we have changed "reduced effect sizes" to "lower signal-to-noise" (revised results section "Resolving functional bases within immune regulatory elements", paragraph 5).

6) Fig 4E: Why would there be an effect outside close interaction range if the two factors influence each other through steric interactions? How did you compute the significance of the counts with a certain distance? What's the null distribution of the selected sites? Why did you not compare to the overall number of selected sites but instead chose the number of increased sizes for a comparison?

7) "Sites that were repressed by BHLHE40 overexpression showed a strong enrichment for motif spacing of 0 to 3 bp, consistent with steric hindrance and competition between factors (FDR < 0.05)..." In contrast to your statements, there is significant repression at 6,8, and 9, whereas there is only significant repression at position 1 and 4. Moreover, there are a significant number of negative effects at various positions and it is unclear whether the shrinking number is due to smaller effects with distance between motifs or because of the distribution of selected motifs. Could you please elaborate on what exactly you mean by the previously reported sites of coordinate GATA Ebox factor binding? Does BHLHE40 have a known activating effect if it is 9bp away from GATA3
Based on reviewer comments and our new data / analyses, we have modified our proposed model of BHLHE40/GATA3 interaction. As discussed above, our data shows opposition of GATA3 and TAL1 to BHLHE40 at RE-4, but does not demonstrate whether this is mediated by steric occlusion between the two factors (see response to major comment 1). We have removed the spacing analysis previously shown in original figure 4E.

7) Methods section: "Model interpretation was conducted as described at <https://github.com/deepmind/deepmind-research/blob/master/enformer/enformer-usage.ipynb>. For CAGE-seq interpretation, we calculated the gradient of the model for unstimulated Jurkat T-cells with respect to the predicted CAGE-seq signal at the CD69 promoter. This was achieved by centering a 393216 bp genomic window within the CD69 promoter (chr12:9760820-9760903)..." Could you briefly describe why you selected regions within the promoter and not downstream or around TSS to model expression?

"... and computing the gradient for human output head # 4831 with respect to output bins 446-450. The absolute value of the gradients were then summed in 128bp bins for coarse grain resolution (Fig 1D).

A similar approach to nominate bases contributing to RE-4 accessibility was adopted to obtain the base resolution contribution scores for the fine-tuned model corresponding to Figure S2 and 3. For this analysis, the window was centered around RE-4(chr12:9764300-9765900) and the gradient was computed with respect to output bins 442-454(Fig S1D, 3D)."

Could you briefly describe why you used more windows than before?

Since CAGE-seq only captures the 5' end of mRNAs, the signal obtained is largely concentrated at the TSS. We therefore followed the recommendations of the Enformer model authors and defined the summed signal at the bins overlapping the annotated TSS as the gene's transcriptional output. These bins correspond to the single bin overlapping the TSS, as well as two bins up/downstream (for a total of five bins). In the case of the fine tuned model where we want to calculate base

contribution scores to the ATAC-seq signal over RE-4, we computed gradients with respect to the predicted model output over the entire RE-4 region(chr12:9764300-9765900), which covered approximately 13 bins of size 128(bins 442-454). We clarify this point in the revised methods.

8) “The implicated interval in RE-4 is over-represented for multiple TF motifs relevant to immune function, including GATA, bHLH/Ebox, TCF, ETS and STAT (Figure 2B).”

What does overrepresentation mean in this context? Did you do a statistical test? Could you describe how you selected (measured over-representation) these five TFs in the Methods section and point to it in the text. If a statistical test was performed, it would be helpful to put a number that describes the overrepresentation in figure 2B. It would also be helpful to put the motifs of these TFs next to their names in Fig2B to help the reader understand what their binding sites have in common or how diverse they are.

We thank the reviewer for their point regarding Fig. 2B. We have removed the phrased “overrepresented”

to reflect that a statistical analysis to determine motif over-representation was not conducted. The TF families highlighted in the figure were selected by first running a standard motif analysis (described in the methods section under “Motif Analysis”), and selecting immune factors implicated in the T-cell activation response. To improve figure clarity, we have added representative motifs for each class of transcription factors to the revised Fig. 2B.

9) “Notably, a second top scoring interval from the CBE and dCas9 screens, centered at sg#48, showed similar TF motif enrichments (Figure 2A; Chr12:9765200-9765310).”

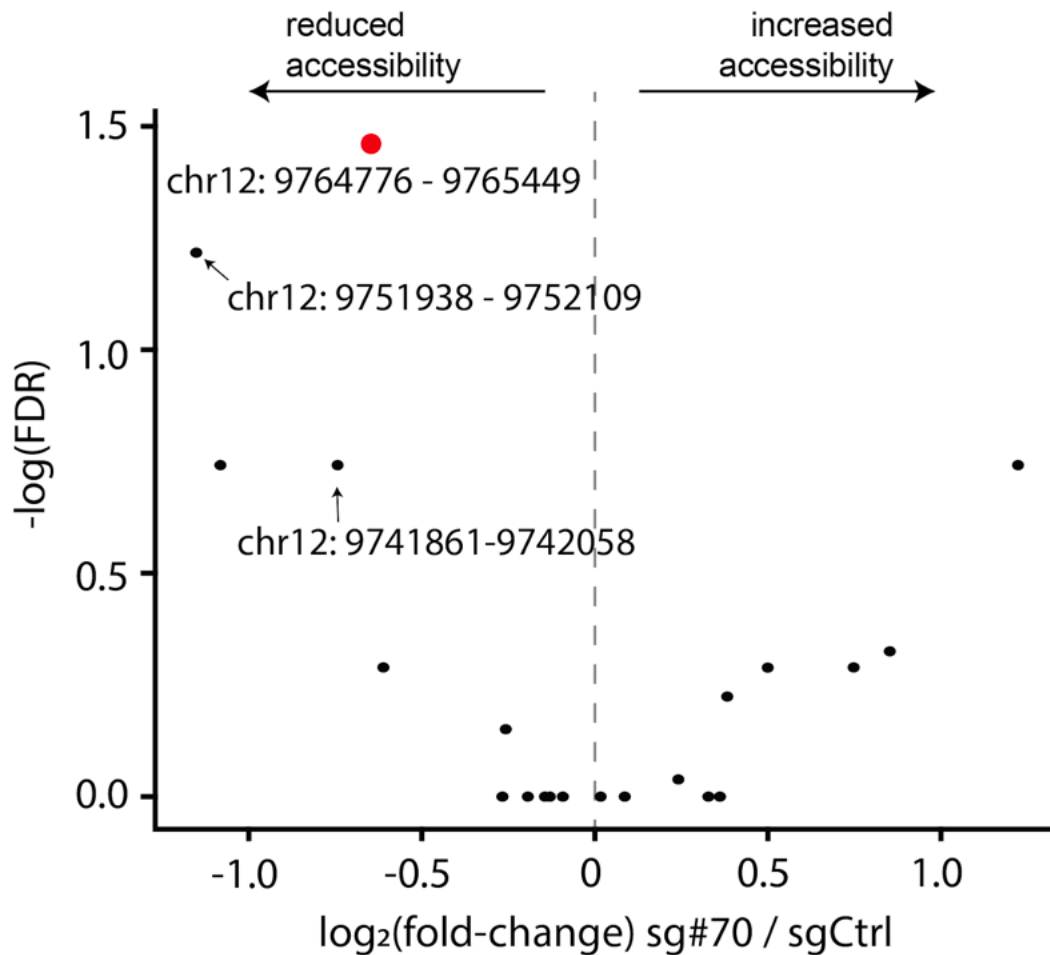
You should at least mention this region earlier in the text because it shows up as early as Fig1D and it is confusing for the reader to wait for this information until now.

We have added a sentence earlier in the results section per reviewer suggestion. We have also added the validation data of sg#48 in revised Fig. S3C.

10) “Fig S3D) Volcano plot depicts chromatin accessibility changes between sgCtrl and sg#70 groups within a 2mb window around RE4(red). X-axis shows $\log_2(\text{fold-change})$ for sg #70 peaks relative to sgCtrl while Y-axis shows $-\log_{10}(\text{FDR})$, with BH correction of p-values based on changes within 1 mb window around RE-4.

Each dot represents changes over how many bases? 1mb seems too large. What’s the dot closest to the red one? Should be explained in more detail in Methods and in the figure caption

Each dot represents an ATAC-seq narrowPeak which ranges in size between ~170 - 1000 bp. Per the reviewer's suggestion, we limited the analysis to peaks within 100 kb of RE-4 (highlighted in red) and labeled the points below. We note that both of these ATAC-seq peaks were targeted with CRISPR-i (revised Fig. 1B-1C) and did not show a significant level of CD69 suppression in comparison to RE-4. We have also added additional details regarding this analysis in the methods and the figure caption.



Reviewer Figure 2.1: Fig. S3E with peaks nearby RE-4 (red) labeled with genomic coordinates

11) "We found that BHLHE40 overexpression suppressed CD69 induction in both control and CBEsg#70 edited Jurkat cells (Figure 4A). However, the magnitude of suppression was greater in the edited cells, potentially due to relief of GATA factor competition." Is the magnitude of suppression really greater? Both are 4-star to ctrl. How to properly compare $90.2 \rightarrow 78.0$ versus $78.2 \rightarrow 54.8$

We have removed this statement regarding the magnitude of suppression

12) "...in our examination of the second interval identified in our dCas9 and CBE screens.

Remarkably, the top base edit hit in this interval (sg#48) also incurs a C->T edit that disrupts a GATA motif flanked by a bHLH/Ebox motif (Figure 2A-2B)"

Could you please provide an enlarged image of the region in the supplementary. Figure 2B shows a lot of GATA/BH combinations at other locations too.

We agree with the reviewer that there are other GATA/BHLH pairs of interest and we now further validate the sg#48 hit. However, we feel that further characterization of this additional region is beyond scope of our study.

13) "We collated all GATA3 bound sites in Jurkat cells that contain a GATA motif and a bHLH/Ebox motif within the corresponding accessible site."

Could you please explain briefly in more detail what was done here or point to the methods section where this is described. How did you collate these GATA3 bound sites in Jurkat cells? I assume you're looking at CHIP-seq from GATA or just motif scans?

The updated analysis uses GATA3 and BHLHE40 bound sites nominated via ChIP-seq for the two

factors in WT stimulated Jurkat cells. It provides evidence that GATA3 and BHLHE40 both regulate a number of immune genes (revised Fig. 4E-4F). We have added further details for the selection of these sites (revised methods section “ChIP-seq processing”) per the reviewer's suggestion.

14) Fig 4G,H

“These results are consistent with a general role of BHLHE40 in restraining GATA3 mediated activation at immune loci. Notably, many of the immune loci subject to opposing regulation contain elements with closely spaced GATA and bHLH/Ebox motifs (0-3 bp), consistent with a general role for competitive TF binding on T cell transcriptional programs and phenotypes (Figure 4G-4H).”

The two presented figures do not possess sufficient resolution to confirm that the close interactions of GATA and BHLHE40 are necessary for the effect. Moreover, there are other positions that show way higher changes in ATAC-accessibility with overexpressed BHLHE40 compared to the regions marked in red. These two figures indicate that BHLHE40 is a global repressor that also represses access of GATA3 to its regulatory element, however it does not confirm that this effect takes place through steric hindrance.

We agree and have revised the model accordingly (see above responses).

Reviewer #2: Chen, et al. present a comprehensive dissection of a regulatory element that controls CD69 expression in stimulated T cells. Their approach involves an innovative combination of regulatory genomics assays, deep neural network predictions, CRISPRi, dCas9 tiling, and CRISPRguided precision base editing. The manuscript convincingly demonstrates that a regulatory element highlighted by differential accessibility and neural network predictions (RE-4) is responsible for CD69 expression in stimulated T cells, and it further shows that GATA and E-box motif elements in RE-4 are strongly contributing to the expression response.

This manuscript represents a new and powerful combination of computational and experimental approaches, where the base-resolution regulatory predictions arising from deep neural networks are tested using base-resolution CRISPR screening techniques. I enjoyed reading the manuscript, but I remain confused and unconvinced by a couple of points.

We appreciate the positive comments and constructive feedback.

1a) The manuscript is a little opaque about exactly what the Enformer neural network contribution scores (i.e., Fig. 1D and Fig S1D) represent and what logic was followed when using the scores to create regulatory predictions. The Enformer scores in Fig 1D appear to be derived from a component of the model that is trained to predict unstimulated T cell CAGE data. Thus, with respect to CD69, the model should be expected to highlight regulatory elements that contribute to the *low* expression of CD69 in unstimulated T cells. It is not clear to me why these contribution scores are used to justify experimental investigation of elements that control higher CD69 expression in stimulated T cells.

The reviewer raises an important point here. We agree that using CAGE data from Jurkat T cells in the stimulated condition would be preferable to nominate regulatory regions using the Enformer model. However, the published model was only trained on resting Jurkat cells, and we could not find public CAGE-seq (or other 5' transcriptomic data for stimulated Jurkat cells in order to fine tune a model. That being said, these Enformer predictions did include most of the genomic sites in the locus that gain accessibility upon stimulation.

We also fine-tuned an Enformer model to predict accessibility in both resting and stimulated Jurkat cells, as well as the differential accessibility between these conditions. Within RE-4, we found that the gradient tracks for the fine tuned model highlighted the same 170 bp interval (revised Fig. 3D, S1E), though with differences in base level predictions (revised Fig. S3G).

1b) According to the Methods, the Enformer model was fine-tuned (i.e., re-trained) to add components that predict ATAC-seq in resting and stimulated T cells. I believe these models produce the contribution score tracks in Fig. S1D (although this should be clarified in the figure caption). However, these models do not seem to be used for any purpose in the investigation of the RE-4 element. Are any sequence elements highlighted as differentially scoring in the stimulated vs. resting

ATAC-seq Enformer models? This would appear to be a more relevant model-driven question than the unstimulated T cell CAGE model scores presented in Fig. 1D.

We appreciate the suggestion. We added an output head to our fine-tuned Enformer model to predict the difference in normalized accessibility between the stimulated and resting conditions, reasoning that the attribution scores from this output head may highlight base pairs that contribute to increased region accessibility during stimulation. The gradient tracks are shown in revised Fig. S1E, 3D, S3G and highlight a similar interval as the gradients from the original model (CAGE-seq output head), as well as a GATA motif and adjacent -e-box/bHLH motifs at the base-pair level. These new data were useful as we refined our model regarding the regulatory contributions of GATA3 and BHLHE40 (see below). We have also added clarification to the revised Fig. S1E regarding finetuned model gradients.

2) Competition between GATA3 and BHLHE40 is proposed to explain expression dynamics in stimulated T cells. This is an appealing model, but I am not fully convinced by the details. Firstly, BHLHE40 is assumed to be a universal repressor, but this is not demonstrated. There are many cases where a TF that is generally repressive can nonetheless play an activating role at a subset of sites. Indeed, BHLHE40 over-expression apparently leads to just as many activated genes as repressed genes in T cells (Fig. 4F). Secondly, the proposed mechanism states that the observed spacing between GATA and E-box motifs is "too close to permit concurrent binding". This is not supported by any modeling or analysis. In the very least, it should be possible for a monomer of BHLHE40 to bind a half-site E-box alongside GATA3; the spacing is not that much less than the observed spacing between GATA motifs and TAL half-site E-boxes in erythroid lineages (PMID: 26503782). Thirdly, the sg#70 CBE edits result in expanded BHLHE40 ChIP-seq signal at RE-4, reduced H3K27ac at RE-4, and reduced expression of CD69 in the context of BHLHE40 overexpression. However, the sg#70 CBE should equally affect both the GATA motif and the neighboring E-box (Fig. 3D). Thus, the sg#70 results cannot by definition support the presented model. Overall, it is not clear why the competition model was favored over activating cooperation between GATA3 and a bHLH TF. It would be interesting to see what effect a BHLHE40 knock-down would have on CD69 expression.

We appreciate these comments, which have prompted extensive computational and experimental analyses for the revision. The new data provide several key insights into TF motifs and binding patterns in the vicinity of the sg70 edited base, allowing us to refine our mechanistic models. First, there are several additional partial e-box motifs in the region, including one site located 9 bp from the GATA motif (revised Fig. 3D, S3G). Our fine-tuned Enformer model highlighted in particular this partial e-box. As intuited by this reviewer, the combination of the GATA motif and this partial e-box completes a potential GATA:TAL1 binding site. We confirmed by ChIP-seq that TAL1 binds the site along with GATA3 in Jurkat cells (revised Fig. 3F). When we next tested binding in edited Jurkat cells, we found that mutation of the GATA motif ablates binding of both TFs, consistent with cooperative TAL1-GATA3 binding to RE-4 in wild-type Jurkat cells.

ChIP-seq also revealed a diffuse increase in BHLHE40 binding across RE-4 and the 5' portion of CD69 in the base-edited Jurkat cells. The diffuse binding suggests the importance of other features, beyond the motif highlighted in our original submission. Indeed, there are several additional (partial) e-box motifs in the region that could contribute to BHLHE40 recruitment. The functional importance of BHLHE40 over the site is further supported by our data showing that BHLHE40 overexpression

decreases RE-4 accessibility and acetylation, and increases RE-4 H3K27 trimethylation (revised Fig. 4C, 4D), and lowers CD69 expression in Jurkat cells (Fig. 4A).

Based on the reviewer points and our new data, we have revised our interpretation of the sg#70 base edit. We simplify our model to state that the base edit displaces the GATA3/TAL1 complex (our sequencing data suggest that the GATA3 motif in specific is the primary target of editing in CBE-sg#70+ cells). The loss of GATA3/TAL1 binding reduces the activity of RE-4 and allows BHLHE40 and potentially other repressive factors to bind the locus. We refer to the possibility

of steric hindrance, but add that this would not be sufficient to explain the relatively widespread changes in factor binding (revised results section “A single artificial variant alters TF binding and suppresses CD69”, paragraphs 5-7; revised Fig. 4E-4F).

Finally, as the reviewer suggested, we did the shRNA KD for BHLHE40 in revised Fig. 4B and confirmed that BHLHE40 KD triggers an induction of CD69 expression. We thank the reviewer for the very helpful comment and hope that the revised and softened model satisfies their concerns.

Reviewer #3: This study by Chen et al. described a method of integrative analysis combining epigenetic perturbations, base editing, and deep learning models for dissecting gene regulatory elements at base resolutions. Taking CD69 gene locus as an example, this analysis identified an artificial C-to-T variant that suppress CD69 expression. This C-to-T base edit was shown to ablate a GATA3 binding site and eliminate GATA3 binding at this site in Jurkat cells. The authors further suggested a potential mechanism of binding competition between GATA3 and BHLHE40 in regulating inducible immune genes and T cell state. While the topic of determining the functions and sequence determinants of cell type-specific regulatory elements is of great interest, this study as presented could not demonstrate that this integrative analysis and machine learning proposed by the authors can help answer this big question.

1) Why does the authors choose to use CRISPRi to test the functional impact of the promoter (RE-3), the putative upstream RE (RE-4) and two other sites (Fig. 1B)? Is this decision purely based on the differential ATAC peak analysis? Did Enformer prediction help here?

In response to reviewer comments, we have generated additional CRISPRi data targeting 3 additional regions nominated by Enformer or differential chromatin accessibility (revised Fig. 1B). These correspond to an intronic region, the gene promoter (as a positive control), and RE-5. These new data enabled us to more systematically compare the CRISPR-i data with the Enformer predictions. Figure S1D shows the relatively good concordance (revised Fig. S1D).

Even though the authors claimed that “the Enformer model predicted that a specific ~170 bp sequence interval within RE-4 is most critical for CD69 regulation”, without this prediction, dCas9 tiling and CBE tiling can also give the same conclusion. The machine learning here did not help narrow down candidates or reduce the amount of screen that need to be performed. After CBE tiling and subsequent sequencing pinpoint C-948 to be critical in regulating CD69 expression, a transcription factor motif analysis can identify its location in a GATA factors binding site. How does the Enformer prediction help in identifying the responsible factor? The machine learning or Enformer prediction does not seem to give new information other than those already provided by other analysis, nor seem to save effort in screening.

We recognize that many conclusions in our paper could have been achieved without the Enformer model. However, our goal here was to evaluate and demonstrate the potential of this model to facilitate regulatory element discovery and dissection. To this effect, the Enformer model was relatively accurate (as we show) and specifically enabled the following:

1. The Enformer model accurately prioritized regulatory elements in the locus, including a larger panel of elements assayed by CRISPRi for the revision (revised Fig. S1D).
2. The model highlighted bases within RE-4 predicted to impact CD69 expression, which in turn led us to identify motifs and nominate cognate TFs (revised Fig. 3D,S3G).
3. The combined analysis provides an early experimental benchmarking of increasingly prevalent deep learning models for genomics (revised Fig. S1D, S2H, S2I).

2) Why is that the authors analyzed BHLHE40, but not BHLHE22 or ARNT2 for their potential role in regulating CD69 expression, while all of them are strongly induced in stimulated Jurkat cells? Even though CBE-sg#70 eliminate GATA3 binding at RE-4, the gain of “broader BHLHE40 binding over RE-4” is not convincing. Could BHLHE22 or ARNT2 show a stronger gain of binding?

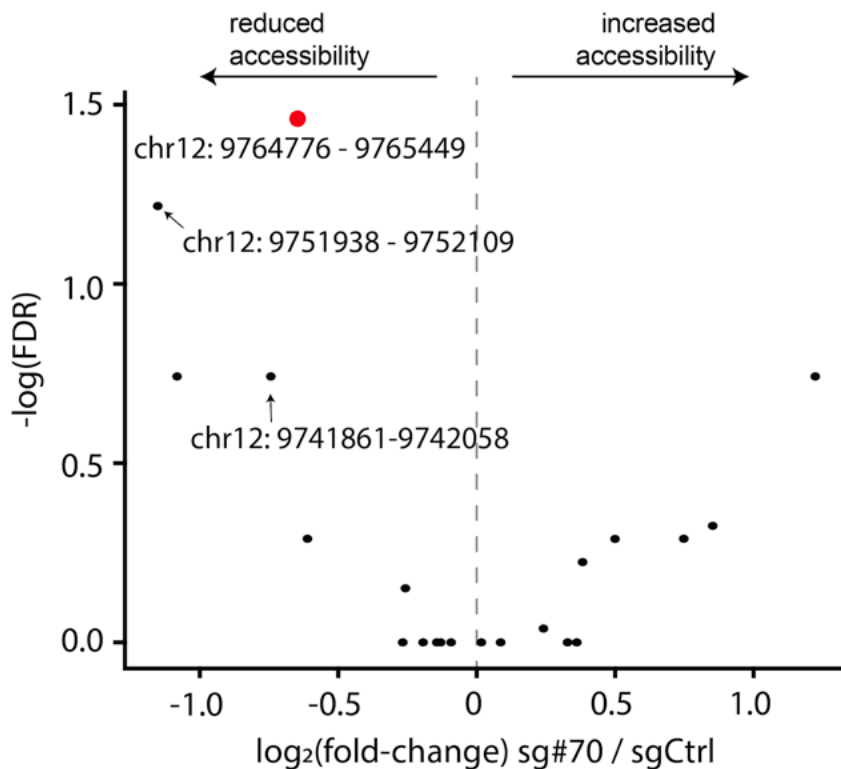
We focused on BHLHE40 due to previous studies implicating the factor in T cell biology, including Th1 function, T cell activation and T cell differentiation⁸⁻¹¹. These papers, though, did not systematically characterize its binding or mechanism of action. We tested the homologous factor BHLHE41, but did not find any evidence that it functions in CD69 regulation (Figure S4D). We

caveat in the revised discussion (paragraph 4) that other factors may also contribute to the regulation of the locus.

3) To support the conclusion that GATA3 and BHLHE40 compete at RE-4 to regulate CD69 expression, could the authors show the effect of BHLHE40 depletion on GATA3 binding at RE-4 and vice versa? How does GATA3-OE and sgBHLHE40 influence CD69 expression?

Based on reviewer comments and our new data, we have revised our interpretation of the sg70 base edit. We simplify our model to state that the base edit displaces the GATA3/TAL1 complex. The loss of GATA3/TAL1 binding reduces the activity of RE-4 and allows BHLHE40 and potentially other repressive factors to bind across the locus. We refer to the possibility of steric hindrance, but add that this would not be sufficient to explain the relatively widespread changes in factor binding (revised results section “A single artificial variant alters TF binding and suppresses CD69”, paragraphs 5-7 ; revised Fig.4E-4F). We do find that GATA3-OE, TAL1-OE or BHLHE40 KD (sgBHLHE40) increase CD69 expression (revised Fig. 4B, S4B-S4C), as is consistent with our model.

4) Could the authors explain what the black dots in Fig. S3D indicates, please? Besides the red dot, a few black ones also show significant change in accessibility. Why does the authors state that "we did not observe any other accessibility changes in the CD69 locus or neighboring genomic regions"? The volcano plot depicts changes in accessibility at ATAC-seq peaks within a 1Mb window centered at RE-4, which corresponds to the red highlighted dot. Only the dots labeled below lie within a 100kb range of the CD69 promoter.



Reviewer Figure 3.1: Fig S3E with peaks nearby RE-4(red) labeled with genomic coordinates We note that both of these ATAC-seq peaks (labeled black dots) were targeted with CRISPRi (Fig. 1B,1C) and did not show a significant level of CD69 suppression in comparison to RE-4 (red dot). We have clarified the legend of Fig. 3E to better explain what each dot represents and the statistical test used. We have also clarified the sentence: “ATAC-seq profiles revealed reduced RE-4 accessibility in cells harboring the CBE-sg#70 construct, relative to CBE controls (Figure 3C,S3D). The effect was most significant for RE-4 (Figure S3E) in the CD69 locus, and we did not observe

significant changes in the vicinity of other activation associated genes such as CD28 and NR4A1 (Figure S3F)." (revised results section "A single artificial variant alters TF binding and suppresses CD69", paragraph 2) .

Reviewer #4: In the manuscript "Integrative dissection of gene regulatory elements at base resolution"

Chen et al. present a framework to identify and validate the impact of single base edits on regulatory element function and target gene expression. For their study they focus on the CD69 gene in T-cells which is induced in T-cell activation. Based on ATAC-seq and machine learning they identify a ~170bp interval within an enhancer predicted to regulate CD69 gene expression. Using base editing they highlight a C-to-T transition that leads to loss of GATA3 binding accompanied by loss of chromatin accessibility. They report that this leads to increased binding of the repressor BHLHE40. They also indicate that this GATA3/BHLHE40 competition is a general genome-wide process during immune cell response and T-cell activation. The current study presents a nice, generalizable experimental and computation framework to characterize the regulatory landscape of an individual locus and insight into potential competition between activators and repressors. The paper is well written and easy to follow with clear experimental and conceptual rationale and execution. Overall, the study provides potentially very interesting insight, but in its current form several aspects of the study seem preliminary. I have several points that would need to be addressed:

We really appreciate the reviewer's positive comments.

1) The advantage of the integration of ATAC and Enformer compared to using the summit or TF footprinting in differential ATAC-seq peaks is not entirely clear to me? Based on both the Enformer and differential accessibility data presented in 1B, would one not predict RE3 to have the strongest effect? What about sites that have strong Enformer signal, e.g. the two peaks in Enformer signal track next to RE2?

This is a good question about our usage of the Enformer model in selecting sites. One aspect of our use of the deep learning model here was to assess Enformer as a tool for prioritizing candidate perturbation sites. As such, sites were selected based on a combination of differential accessibility between the resting/stimulated conditions in both Jurkat and CD4 T-cells and the Enformer gradient score (see revised methods). In the revised manuscript we have included data showing CRISPR-i targeting of 3 additional regions(revised Fig 1B), which correspond to an intronic region, the gene promoter (as a positive control), and RE-5. We have also included a comparison of the CRISPR-i data with the predicted impact of each guide based on the Enformer gradient score and found generally good agreement between the two (revised Fig. S1D).

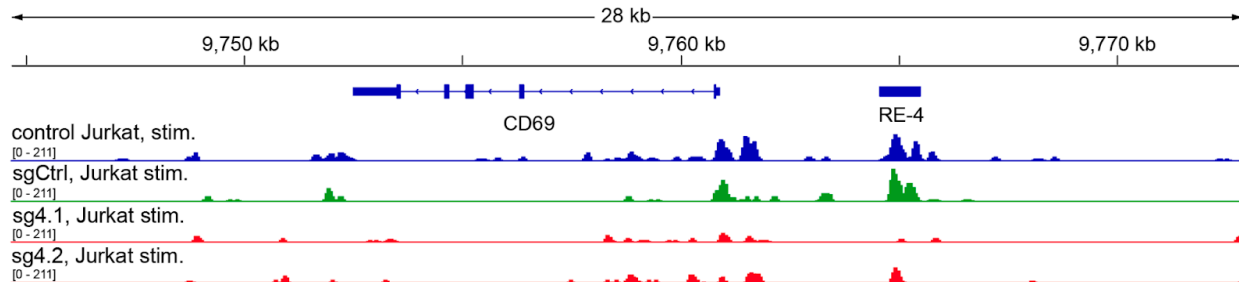
RE-3 lies approximately ~2-2.5kb away from both RE-4 and the CD69 promoter. Given the expected perturbation size of dCas9-KRAB(~1-2 kb), we would expect a strong effect of CRISPR-i perturbation of this region. Based on the Enformer gradient score, we predict that perturbation of RE-3 would show the third strongest effect in our set of targeted regions, after the promoter and RE-4. As seen in Fig. S1D, this prediction matches well with the experimental results.

The two gradient peaks proximal to RE-2 referenced by the reviewer were not included in the CRISPR-i library, since one overlapped an exon/3'UTR of the gene, while the other did not overlap an accessible site in either resting/stimulated Jurkats/CD4 T-cells.

2) Since the promoter is required for gene expression and indeed are often used as positive controls in screens it is surprising to see only mild effect when targeting the promoter proximal RE3. In addition to focusing on CD69 protein levels, ATAC-seq data (or H3K27ac) would be helpful to show that chromatin accessibility and activity of the respective REs is lower after infection with guide RNAs. This would provide a more direct readout and enable a mechanistic link that indeed RE activity is linked to expression.

Though the RE-3 is promoter proximal as the reviewer indicates, it still lies approximately more than 2.5kb away from the gene promoter (see response to major comment 1 for further discussion). Given that dCas9-KRAB impacts a region size of ~1-2kb, we do not expect (or predict, using

Enformer gradient) that an RE-3 perturbation would have as strong an effect size as targeting the gene promoter. We have now added data corresponding to direct CRISPR-i targeting of the gene promoter, which (as expected) has a stronger effect on CD69 expression (revised Fig. S1C) than targeting RE-3. Also, per reviewer suggestion, we provide ATAC-seq data confirming significantly reduced accessibility over CRISPR-i targeted sites.



Reviewer Figure 4.1: Plot of CD69 locus demonstrating reduced ATAC-seq signal in Jurkats with dCas9-KRAB and RE-4 targeting guides sg4.1, sg4.2 (red), relative to a non-targeting control (sgCtrl, green). Control, stimulated Jurkats without the dCas9-KRAB construct (blue) are also shown. ATAC-seq tracks are RPGC normalized to 1x coverage.

3) Could the modest effect of sg#70 and overexpression of BHLH40 on chromatin accessibility, e.g. Fig. 4B and C might be because silencers are also accessible? H3K27ac or H3K27me3 levels might be better chromatin marks to study the competition of activator and repressor. For n=2 experiments please display individual values. The statistics using uncorrected p-value despite having whole genome data with thousands of peaks seems strange, particularly since there are many more pronounced sites as shown in panel D of the same figure.

The reviewer raises an excellent point regarding the potential role of silencers in regulating CD69. We gathered H3K27ac and H3K27me3 ChIP-seq data in stimulated BHLH40 overexpressing(OE) Jurkat cells (revised Fig. 4D). We observe that BHLH40OE leads to increased BHLH40 binding over RE-4, decreased H3K27ac and increased H3K27me3 signal over RE-4. These data are consistent with a more repressive state within RE-4 during BHLH40OE.

We have added in the individual values for figure 4B(now figure 4C). We have also changed the uncorrected p-value referenced in the figure legend to a genome-wide FDR.

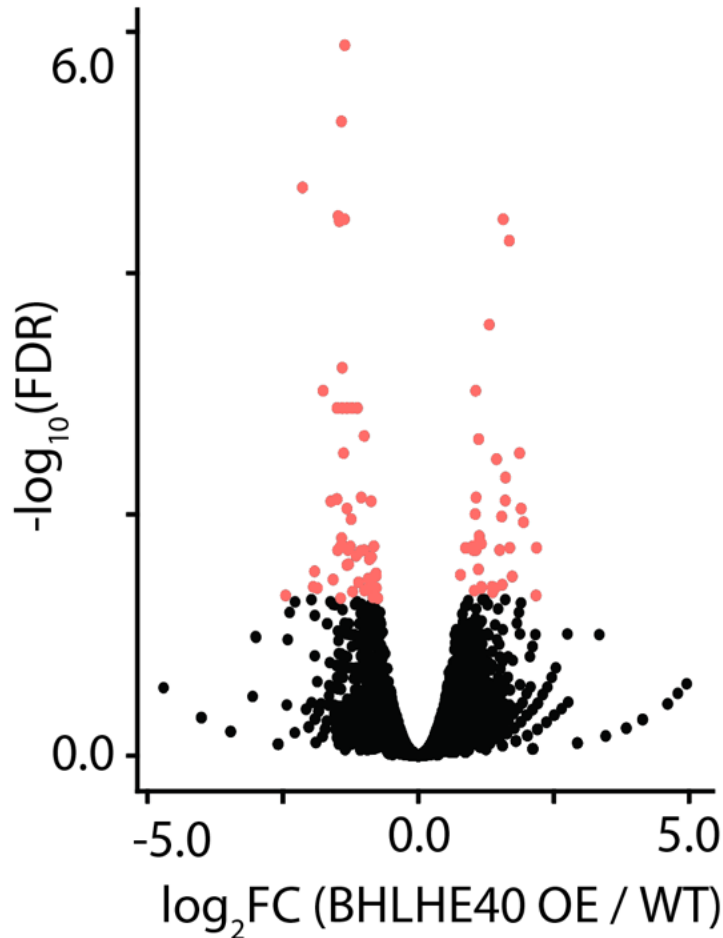
4) Sg#70 seems to have the most drastic effect on recruitment of BHLH40 to the promoter or proximal RE of CD69, whereas the signal at RE4 is comparable and BHLH40 is already bound in the WT (Fig.3F). If GATA3 as stated prevents binding of BHLH40 why is there no clear difference in binding at RE4? How is the binding pattern of BHLH40 upon infection with Sg#70 and overexpression of BHLH40?

The reviewer makes a good point as to the nature of exclusive binding between GATA3 and BHLH40 at the site of the base edit. Based on multiple reviewer comments and our new data, we have revised our interpretation of the sg70 base edit. We simplify our model to state that the base edit displaces the GATA3/TAL1 complex. The loss of GATA3/TAL1 binding reduces the activity of RE-4 and allows BHLH40 and potentially other repressive factors to bind across the locus. We refer to the possibility of steric hindrance, but add that this would not be sufficient to explain the relatively widespread changes in factor binding (results section "A single artificial variant alters TF binding and suppresses CD69", paragraphs 5-7; revised Fig.4E-4F).

5) BHLH40 is described as transcriptional repressor, but Fig 4F illustrates up and downregulated genes linked to BHLH40 REs. How do the authors explain this? Can BHLH40 context specific activate or repress? How were target genes defined? Is this a mix of direct and indirect effects? Could these be resolved?

BHLH40 has been primarily described as a transcriptional repressor in prior studies^{10,12-16}. Overall, BHLH40 over-expression leads to a global decrease in accessibility (Revised Fig. S5A). However, BHLH40 over-expression leads to similar numbers of differentially up and down

regulated genes (revised Fig. S5B). We also analyzed how increased BHLHE40 binding impacts chromatin accessibility. We collated the set of BHLHE40 ChIP-seq peaks in WT Jurkat cells and subsetted to those that gain BHLHE40 binding signal in the BHLHE40 OE condition. These sites showed roughly equal numbers of differentially accessible/repressed ATAC-seq peaks (reviewer Figure 4.2).



Reviewer Figure 4.2: Volcano plot of ATAC-seq signal change between BHLHE40 OE and WT(Ctrl-LV) over BHLHE40 ChIP-seq peaks(see revised Methods) that gain binding signal during BHLHE40OE. Red points represent differentially accessible ATAC-seq peaks between the two conditions at FDR = 0.05.

In order to resolve whether the gene expression changes are direct or indirect effects, we took all DE genes in BHLHE40OE and determined whether they had a TSS proximal (< 25kb) BHLHE40 ChIP-seq peak in WT Jurkat cells. We found that the majority (739/1007) DE genes did have TSS proximal BHLHE40 binding (Revised Fig. S5D-S5E).

Although BHLHE40 may have a greater tendency to act as a repressor, we agree with the reviewer that based on these results that it is likely a context specific regulator that can repress as well as activate certain immune genes (Revised Fig. S5C,S5E,S5F). We have modified our statements in the main text accordingly (revised results section “A single artificial variant alters TF binding and suppresses CD69”, paragraph 6 and “Global interplay between BHLHE40 and GATA3 in T cell response”, paragraphs 1-2).

6) The authors propose that GATA3 and BHLHE40 competition impacts transcriptional responses globally. I think to show this it would be important not only to show changes after overexpression of BHLHE40 since this can lead to artifacts. It would be better to integrate RE activity, binding of

GATA3 and BHLHE40 by ChIP-seq and RNA-seq (several of these datasets are already generated as part of this study) in control and stimulated T-cells. This could directly address several questions of competition, for example for genes that are higher expressed in activated T-cells do they gain GATA3, lose BHLHE40 at REs or does the relative ratio of binding change? How many sites are only bound by GATA3 or BHLHE40 and how many are co-bound? What about the associated H3K27ac and accessibility levels? And similar analysis could be done for downregulated genes. How many genes are dependent on the competition in this model?

We appreciate the comment. Based on this comment and others, we extended our analysis of global interactions between GATA3 and BHLHE40 within stimulated Jurkat T-cells by integrating ChIP-seq data for BHLHE40, GATA3, and H3K27ac within stimulated WT Jurkat cells (revised results section “Global interplay between BHLHE40 and GATA3 in T-cell responses;”, methods section “ChIP-seq processing”; revised Fig. 4E-F, S5A-F).

First, we found that BHLHE40OE altered the expression of multiple immune targets, with a general up-regulation of Th1 and effector T-cell related genes and down-regulation of Th2 and naive/stemness pathways (Fig. S5B-S5C). These pathway changes oppose those reported for GATA3, which has been implicated in Th2 differentiation and naive/stemness phenotypes.

We then defined active enhancers bound by either BHLHE40 or GATA3 by intersecting either factor’s IDR ChIP-seq peaks with H3K27ac peaks. We found that a large number of enhancers were bound by both BHLHE40 and GATA3 (10623 / 21146 candidate BHLHE40 bound enhancers; Fig. 4E). We also collated the set of associated genes for each factor (defined by the presence of binding event < 25kb from an annotated TSS), and found that the two factors share a large number of target genes (5479; Fig. 4E). Generally, this set of BHLHE40/GATA3 co-regulated genes showed expression changes during BHLHE40OE consistent with up-regulation of effector T-cell related pathways (Fig. 4F).

These data support a general interaction between the GATA3 and BHLHE40 at key immune loci.

Rebuttal references:

1. Sanda, T., Lawton, L.N., Barrasa, M.I., Fan, Z.P., Kohlhammer, H., Gutierrez, A., Ma, W., Tatarek, J., Ahn, Y., Kelliher, M.A., et al. (2012). Core transcriptional regulatory circuit controlled by the TAL1 complex in human T cell acute lymphoblastic leukemia. *Cancer Cell* 22, 209–221. 10.1016/j.ccr.2012.06.007.
2. De Masi, F., Grove, C.A., Vedenko, A., Alibés, A., Gisselbrecht, S.S., Serrano, L., Bulyk, M.L., and Walhout, A.J.M. (2011). Using a structural and logics systems approach to infer bHLH–DNA binding specificity determinants. *Nucleic Acids Res.* 39, 4553–4563. 10.1093/nar/gkr070.
3. Grant, C.E., and Bailey, T.L. XSTREME: Comprehensive motif analysis of biological sequence datasets. 10.1101/2021.09.02.458722.
4. Rees, H.A., and Liu, D.R. (2018). Base editing: precision chemistry on the genome and transcriptome of living cells. *Nat. Rev. Genet.* 19, 770–788. 10.1038/s41576-018-0059-1.
5. Hanna, R.E., Hegde, M., Fagre, C.R., DeWeirdt, P.C., Sangree, A.K., Szegletes, Z., Griffith, A., Feeley, M.N., Sanson, K.R., Baidi, Y., et al. (2021). Massively parallel assessment of human variants with base editor screens. *Cell* 184, 1064–1080.e20. 10.1016/j.cell.2021.01.012.
6. Shiraki, T., Kondo, S., Katayama, S., Waki, K., Kasukawa, T., Kawaji, H., Kodzius, R., Watahiki, A., Nakamura, M., Arakawa, T., et al. (2003). Cap analysis gene expression for high-throughput analysis of transcriptional starting point and identification of promoter usage. *Proc. Natl. Acad. Sci. U. S. A.* 100, 15776–15781. 10.1073/pnas.2136655100.
7. Avsec, Ž., Agarwal, V., Visentin, D., Ledsam, J.R., Grabska-Barwinska, A., Taylor, K.R., Assael, Y., Jumper, J., Kohli, P., and Kelley, D.R. (2021). Effective gene expression prediction from sequence by integrating long-range interactions. *Nat. Methods* 18, 1196–1203. 10.1038/s41592-021-01252-x.
8. Zhang, L., Yu, X., Zheng, L., Zhang, Y., Li, Y., Fang, Q., Gao, R., Kang, B., Zhang, Q., Huang, J.Y., et al. (2018). Lineage tracking reveals dynamic relationships of T cells in colorectal cancer. *Nature* 564, 268–272. 10.1038/s41586-018-0694-x.

9. Lin, C.-C., Bradstreet, T.R., Schwarzkopf, E.A., Sim, J., Carrero, J.A., Chou, C., Cook, L.E., Egawa, T., Taneja, R., Murphy, T.L., et al. (2014). Bhlhe40 controls cytokine production by T cells and is essential for pathogenicity in autoimmune neuroinflammation. *Nat. Commun.* 5, 3551. 10.1038/ncomms4551.
 10. Cook, M.E., Jarjour, N.N., Lin, C.-C., and Edelson, B.T. (2020). Transcription Factor Bhlhe40 in Immunity and Autoimmunity. *Trends Immunol.* 41, 1023–1036. 10.1016/j.it.2020.09.002.
 11. Yu, F., Sharma, S., Jankovic, D., Gurram, R.K., Su, P., Hu, G., Li, R., Rieder, S., Zhao, K., Sun, B., et al. (2018). The transcription factor Bhlhe40 is a switch of inflammatory versus antiinflammatory Th1 cell fate determination. *J. Exp. Med.* 215, 1813–1821. 10.1084/jem.20170155.
 12. Huynh, J.P., Lin, C.-C., Kimmey, J.M., Jarjour, N.N., Schwarzkopf, E.A., Bradstreet, T.R., Shchukina, I., Shpynov, O., Weaver, C.T., Taneja, R., et al. (2018). Bhlhe40 is an essential repressor of IL-10 during Mycobacterium tuberculosis infection. *J. Exp. Med.* 215, 1823–1838. 10.1084/jem.20171704.
 13. Emming, S., Bianchi, N., Polletti, S., Balestrieri, C., Leoni, C., Montagner, S., Chirichella, M., Delaleu, N., Natoli, G., and Monticelli, S. (2020). A molecular network regulating the proinflammatory phenotype of human memory T lymphocytes. *Nat. Immunol.* 21, 388–399. 10.1038/s41590-020-0622-8.
 14. Asanoma, K., Liu, G., Yamane, T., Miyanari, Y., Takao, T., Yagi, H., Ohgami, T., Ichinoe, A., Sonoda, K., Wake, N., et al. (2015). Regulation of the Mechanism of TWIST1 Transcription by BHLHE40 and BHLHE41 in Cancer Cells. *Mol. Cell. Biol.* 35, 4096–4109. 10.1128/MCB.00678-15.
 15. Zawel, Yu, and Torrance DEC1 is a downstream target of TGF- β with sequence-specific transcriptional repressor activities. *Proc. Estonian Acad. Sci. Biol. Ecol.*
 16. Honma, S., Kawamoto, T., Takagi, Y., Fujimoto, K., Sato, F., Noshiro, M., Kato, Y., and Honma, K.-I. (2002). Dec1 and Dec2 are regulators of the mammalian molecular clock. *Nature* 419, 841–844. 10.1038/nature01123.
-

Referees' report, second round of review

Reviewer #1: Comments enter in this field will be shared with the author; your identity will remain anonymous.

Reviewer #2: The authors have satisfactorily addressed my previous comments with additional analyses and experiments, and by revising some of their claims.

Reviewer #3: The authors generated additional data targeting 3 additional regions, in which the intronic region and gene promoter region were nominated by both Enformer and differential chromatin accessibility, RE-5 nominated purely by differential chromatin accessibility. The authors did not directly explain how Enformer prediction helps with identifying candidate regulatory sites other than confirming differential ATAC-seq peak analysis. The Enformer prediction model is trained on chromatin maps. The authors have shown that Enformer prediction was relatively accurate in prioritizing regulatory elements, but Enformer prediction did not provide additional data that a direct ATAC-seq differential peak analysis could not provide. The Enformer model relies on all other ATAC-seq, CRISPR-i, dCas9 and base editing data to perform prediction, it may show greater contribution in the future, but as it currently stands, this deep learning model had not shown clear benefits in identifying gene regulatory elements. This study combined ATAC-seq, CRISPR-i, dCas9 and base editing analysis in identifying RE-4, C-948, GATA-3, TAL-1 and BHLHE40 in regulating CD69 expression, but all the methods used here are already well-established, the sites and factors identified are also partially required for CD69 expression.

Overall the priority and novelty of the manuscript does not meet the level for Cell Genomics, and may be more suitable for another journal.

Reviewer #4: The authors have done a great job addressing the comments.

Minor point:

... while reducing accessibility and H3K27ac over the element (Figure 4C).

Please change since H3K27ac is displayed in panel 4D. Based on the displayed tracks the change in H3K27ac signal at RE-4 is not visible to me.

Authors' response to the second round of review

Response to reviewer comments 3/23/23

Summary of additions and revisions

We have added a separate discussion section on the limitations of our study, and specifically address limitations related to the use of the deep learning model that reviewer 3 raised. We have also fixed panel labels in Figure 4 in the main text, related to the error noted by reviewer 4. We additionally have reduced word counts and reformatted parts of the methods section where necessary.