

Integrative approaches generate insights into the architecture of non-syndromic cleft lip with or without cleft palate

Julia Welzenbach,¹ Nigel L. Hammond,² Miloš Nikolić,³ Frederic Thieme,¹ Nina Ishorst,¹ Elizabeth J. Leslie,⁴ Seth M. Weinberg,^{5,6} Terri H. Beaty,⁷ Mary L. Marazita,^{5,6,8} Elisabeth Mangold,¹ Michael Knapp,⁹ Justin Cotney,^{10,11} Alvaro Rada-Iglesias,^{3,12,13} Michael J. Dixon,² and Kerstin U. Ludwig^{1,*}

Summary

Non-syndromic cleft lip with or without cleft palate (nsCL/P) is a common congenital facial malformation with a multifactorial etiology. Genome-wide association studies (GWASs) have identified multiple genetic risk loci. However, functional interpretation of these loci is hampered by the underrepresentation in public resources of systematic functional maps representative of human embryonic facial development. To generate novel insights into the etiology of nsCL/P, we leveraged published GWAS data on nsCL/P as well as available chromatin modification and expression data on mid-facial development. Our analyses identified five novel risk loci, prioritized candidate target genes within associated regions, and highlighted distinct pathways. Furthermore, the results suggest the presence of distinct regulatory effects of nsCL/P risk variants throughout mid-facial development and shed light on its regulatory architecture. Our integrated data provide a platform to advance hypothesis-driven molecular investigations of nsCL/P and other human facial defects.

Introduction

Current research into the etiology of common disorders is focused on the identification of genetic susceptibility factors and the manner in which these risk variants interfere with biological function. Over the past decade, genome-wide association studies (GWASs) of common disorders have identified numerous risk loci. However, success in the translation of statistical associations from GWASs into functional mechanisms is only a very recent achievement.^{1–6} A major driver of these advances has been the availability of large-scale genetic data and the systematic integration of genetic, transcriptional, epigenetic, and other -omics datasets from disease-relevant cell types and tissues.⁷

Facial disorders rank among the most common birth defects worldwide and represent a substantial burden for affected individuals, their families, and healthcare systems.^{8,9} The most frequent facial disorder is non-syndromic cleft lip with or without cleft palate (nsCL/P). This condition has a global incidence of ~1 in 1,000 live births⁹ and is characterized by a multifactorial etiology that includes an overall genetic contribution of around

90%.^{9–11} On an epidemiological level, nsCL/P is associated with an increased risk for adverse health outcomes.¹² However, this observation remains largely unexplained at both the clinical and molecular levels. To date, GWASs and other systematic approaches have identified at least 40 nsCL/P risk loci,^{13–28} which explain up to 30% of the estimated heritability in European populations.²¹ Despite these successes, functional dissection of the associated regions has been limited to only a few loci.^{29–32} This is mainly attributable to the systematic underrepresentation of embryonic facial data in public resources such as ENCODE,³³ Roadmap Epigenome,³⁴ and GTEx.³⁵ To overcome this limitation, researchers have recently profiled multiple chromatin modifications in cell types and tissues of relevance to individual time points of mid-facial development, a process that is largely completed by week 10 of gestation (Figure 1A). These cell types and tissues include early human neural crest cells (hNCCs),³⁷ lineage-specified human cranial NCCs (cNCCs),³⁸ and embryonic mid-facial tissue samples encompassing the time period 4.5–10 weeks post-conception (craniofacial tissue [CT]; days 32–56 of gestation).³⁹ Previous studies have demonstrated a significant enrichment of nsCL/P-GWAS

¹Institute of Human Genetics, University Hospital Bonn, Medical Faculty, Venusberg-Campus 1, 53127 Bonn, Germany; ²Faculty of Biology, Medicine, and Health, Manchester Academic Health Sciences Centre, University of Manchester, Manchester M13 9PT, UK; ³Center for Molecular Medicine Cologne (CMCC), University of Cologne, Cologne, Germany; ⁴Department of Human Genetics, School of Medicine, Emory University, Atlanta, GA, USA; ⁵Center for Craniofacial and Dental Genetics, Department of Oral Biology, School of Dental Medicine, University of Pittsburgh, Pittsburgh, PA, USA; ⁶Department of Human Genetics, Graduate School of Public Health, University of Pittsburgh, Pittsburgh, PA, USA; ⁷Department of Epidemiology, Bloomberg School of Public Health, Johns Hopkins University, Baltimore, MD, USA; ⁸Department of Psychiatry and Clinical and Translational Science Institute, School of Medicine, University of Pittsburgh, Pittsburgh, PA 15260, USA; ⁹Institute of Medical Biometry, Informatics, and Epidemiology, University Hospital Bonn, Bonn, Germany; ¹⁰Department of Genetics and Genome Sciences, UConn Health, Farmington, CT, USA; ¹¹Institute for Systems Genomics, University of Connecticut, Storrs, CT, USA; ¹²Cluster of Excellence Cellular Stress Responses in Aging-Associated Diseases (CECAD), University of Cologne, Cologne, Germany; ¹³Institute of Biomedicine and Biotechnology of Cantabria (IBBT), University of Cantabria, Cantabria, Spain

*Correspondence: kerstin.ludwig@uni-bonn.de

<https://doi.org/10.1016/j.xhgg.2021.100038>

© 2021 The Authors. This is an open access article under the CC BY-NC-ND license (<http://creativecommons.org/licenses/by-nc-nd/4.0/>).



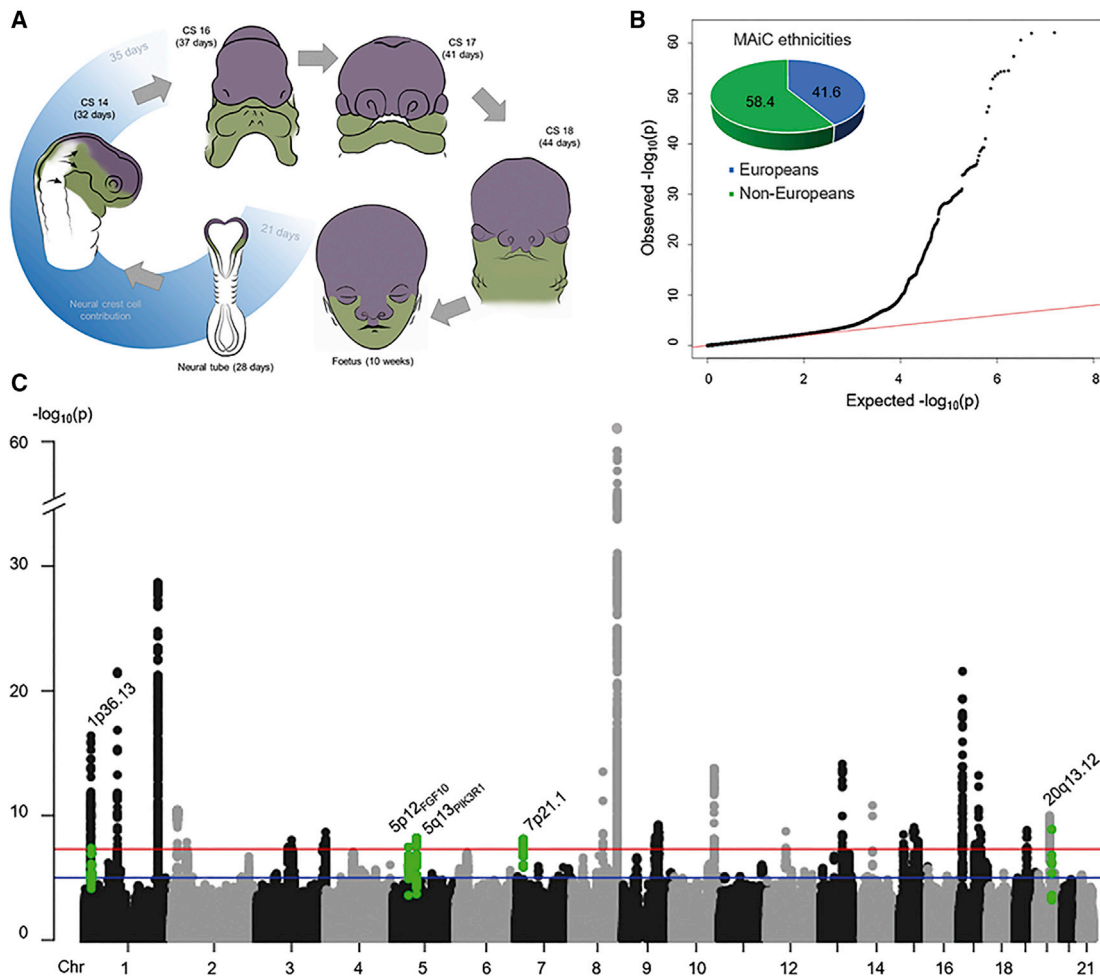


Figure 1. Human facial development and results of meta-analysis in clefting (MAiC)

(A) Schematic representation. The first phase of facial development (blue shading) is characterized by a substantial contribution of neural crest cells (NCCs): In early embryogenesis, NCCs arise in the ectoderm, undergo epithelial-to-mesenchymal transition, and begin to migrate from the dorsal neural tube. An NCC fraction (i.e., cranial NCCs) contribute to the pre-swellings of the face and populate the future frontonasal prominence as well as the first (purple) and second (green) pharyngeal arches.³⁶ Subsequently, NCC-derived cells fuse to form those human facial structures that are finalized by the 10th week of embryogenesis.

(B) MAiC quantile-quantile plot. Observed statistical associations for non-syndromic cleft lip with/without cleft palate (nsCL/P) are plotted against the association statistics expected under the null hypothesis of no association. The contribution of different ethnicities in MAiC is shown using a pie chart.

(C) MAiC Manhattan plot. MAiC $-\log_{10}(p)$ association results are plotted along their chromosomal distribution. Blue and red lines indicate suggestive ($p < 10^{-5}$) and genome-wide ($p < 5 \times 10^{-8}$) significance, respectively. The lowest p value was observed for rs55658222 ($p = 8.69 \times 10^{-63}$), located at 8q24.²⁷ Novel risk loci are highlighted in green (lead variant plus variants in linkage disequilibrium [LD] [$r^2 \geq 0.6$]). Gene names in subscript discriminate novel risk loci in situations where the respective chromosomal band is already listed among the 40 risk loci.

variants in active chromatin regions from both hNCCs and CT.^{21,39} To date, however, the fact that these datasets have been generated from differing sources has precluded the integrative analyses required for a comprehensive assessment of variant function at different time points of mid-facial development.

To generate novel insights into the etiology of nsCL/P, the present study leveraged both existing GWAS data on nsCL/P and epigenetic data on mid-facial development. The specific aims of the study were threefold (Figure S1). First, we generated one of the largest genome-wide genetic datasets for nsCL/P to date by combining three GWASs, which collectively encompassed European, Asian, and

Latin American ethnicities. Using this resource, which we term MAiC (meta-analysis in clefting), we confirmed the vast majority of established risk regions and detected five novel loci (the strategy for identification of novel risk loci is described in the Supplemental Material and methods). To shed light on potential etiological overlaps between nsCL/P and other phenotypes, we then cross-referenced the lead variants at nsCL/P risk loci with GWAS data on >3,000 common traits and identified a set of loci with pleiotropic effects. Second, we compiled a comprehensive epigenetic map of mid-facial development through joint analyses of available data from hNCCs, cNCCs, and CT. This resource of chromatin segments

across mid-facial development serves as a platform for the interpretation of genetic findings for facial disorders and traits. Finally, we aimed to generate systematic insights into nsCL/P biology by combining MAiC and epigenetic data and then adding additional layers on gene expression in NCCs and global and local three-dimensional (3D) genomic interactions (i.e., topologically associated domains [TADs],⁴⁰ promoter-capture HiC [pChI-C]⁴¹). This approach revealed tissue- and time-point-specific regulatory effects at GWAS risk loci, prioritized candidate target genes, and highlighted distinct pathways. To our knowledge, the present report is the first to describe the systematic integration of large-scale summary statistics in nsCL/P and data on the *cis*-regulatory landscape across several stages of human mid-facial development.

Material and methods

GWAS meta-analysis MAiC

Cohort description

The meta-analysis included data from three previously published individual GWASs on nsCL/P (Bonn case-control GWAS cohort,¹⁸ GENEVA trio cohort,²⁰ POFC GWAS cohort;¹⁷ Table S1). We included all nsCL/P summary statistics that were publicly accessible until June 2018. Data from the Bonn cohort were available in-house, while both the GENEVA (dbGaP: phs000094) and POFC (dbGaP: phs000774) datasets were downloaded from dbGaP upon approved data access, respectively. Previously conducted meta-analyses included combinations of two of these studies (Bonn and GENEVA GWAS cohort in Ludwig et al., 2012¹⁹ [genotyped variants] and 2017²¹ [imputed variants], GENEVA and POFC in Leslie et al.²⁶). In the present study we combined the three GWAS cohorts to generate the largest nsCL/P meta-analysis to date. In accordance with previous studies,^{19,21,26} two meta-analyses were performed: (1) using all individuals with diverse population backgrounds (to increase statistical power by maximizing sample size; in the following termed as MAiC), and (2) using the European datasets only (MAiC_{Euro}, to reduce genetic heterogeneity based on population differences). Data quality control (QC) included the detection and removal of overlapping individuals, confirmation of ethnicity, and data re-analysis. We call this new dataset MAiC to provide a clear distinction from the previous individual studies and meta-analyses of sub-cohorts. Further details in cohort description and data QC can be found in the [Supplemental information](#).

Statistical analyses

Statistical analyses were performed separately for case-control cohorts and case-parent trios, respectively. Imputed data were taken as provided by dbGaP (POFC) or generated as previously described (for Bonn and GENEVA),²¹ respectively, and best-guess genotypes were assigned based on *a posteriori* genotype probabilities of ≥ 0.6 . In the case-control cohorts, GWAS was performed using logistic regression performed with SNPTEST and -method expected, by incorporating five (Bonn and GENEVA cohorts) and 18 (POFC cohort) dimensions of the multi-dimensional-scaling coordinates,⁴² respectively. For the case-parent trios, a transmission disequilibrium test (TDT) was performed on the best-guess genotypes.⁴³ After data cleaning procedures ([Supplemental information](#)), we meta-analyzed the GWAS data of all four sub-cohorts (Bonn case-control, GENEVA case-parent trios, POFC case-control, and POFC case-parent trios) using METAL.⁴⁴

The final MAiC dataset (case-control plus case-parent trios) contained 6,825 individuals (including 3,946 affected; MAiC_{Euro}: 3,568 individuals including 1,517 affected; Table S1). The maximum genomic inflation factor was 1.051 (GENEVA) and 1.056 (POFC case-control) for MAiC and MAiC_{Euro}, respectively. All functional downstream analyses are based on MAiC because of largely increased statistical power. To estimate the single-nucleotide polymorphism (SNP)-based heritability (h^2) for nsCL/P on the liability scale, we generated a European case-control-only dataset (Bonn, POFC, totaling 532 cases and 2,051 controls; Table S1) and performed linkage disequilibrium (LD) score regression as implemented in ldsr.⁴⁵ Sample and population prevalence were set to 0.21 and 0.001, respectively.

Gene-based and pathway analyses

Gene-based analyses in MAiC and MAiC_{Euro} were performed using MAGMA⁴⁶ (v.1.06), implemented in FUMA. The input SNPs of MAiC were mapped to 17,911 protein-coding genes based to a distance of 0 kb upstream/downstream of the genes, resulting in threshold of test-wide significance of $p = 2.79 \times 10^{-6}$ (i.e., 0.05/17,911). To annotate known and novel nsCL/P risk loci in biological context, we investigated common expression patterns of the GWAS_{TAD} genes and their molecular functions (gene ontology [GO] terms, Kyoto Encyclopedia of Genes and Genomes (KEGG) pathways) using FUMAs "GENE2FUNC" tool in (1) all GWAS_{TAD} genes, and (2) a subset of GWAS_{TAD} genes expressed in NCCs. This approach allows us to pinpoint risk loci or genes that are functionally involved in the same pathways or molecular processes and might be useful for gene prioritization.

Analysis of pleiotropic effects using the GWAS ATLAS

For each of the 45 lead SNPs in MAiC, association signals from large-scale genetic studies (including p value, effect size, and effect direction) were retrieved from the GWAS ATLAS.⁴⁷ At time of analysis (November 2019), the database comprised 4,756 GWASs on 3,302 unique traits. Notably, the unique traits are split into 28 domains, of which we combined two (environment, activities) into one domain to reduce redundancy. All significant SNP-trait associations at $p < 0.05$ were considered, and this number was corrected for the number of GWASs and loci in the analysis.

Epigenetic datasets for mid-facial development

Identification of datasets relevant to mid-facial development

Human cell-type- and developmental-stage-specific data for mid-facial development are underrepresented (or not represented at all) in large consortia data such as ENCODE.³³ However, available data in the Gene Expression Omnibus (GEO) covered mid-facial development from (1) early stages (hNCCs,³⁷ accessed through GEO: GSE28874), (2) differentiated human cNCCs³⁸ (accessed through GEO: GSE70751), and (3) embryonic craniofacial human tissue of different Carnegie stages (CS) (accessed through GEO: GSE97752).³⁹ In each of these datasets, analyses of chromatin modifications were performed using chromatin immunoprecipitation followed by sequencing (chromatin immunoprecipitation sequencing [ChIP-seq]) or are available as imputed datasets. Detailed information including antibodies used in these studies is shown in Table S3 and in the [Supplemental information](#). For hNCCs and cNCCs, ChIP-seq had been performed for chromatin modifications H3K27ac, H3K4me1, H3K4me3, and H3K27me3. In CT, for samples of CS13–CS17, ChIP-seq was performed for H3K27ac, H3K4me1, H3K4me3, H3K27me3, and H3K36me3 (Table S4), and data for H3K9me3 were imputed. For CS20 and

10 wpc, H3K27ac3 ChIP-seq data were experimentally derived; all other marks were imputed (Table S3).

Data processing

For hNCCs and cNCCs, raw data were available in fastq format. A description of data QC is given in Rada-Iglesias et al.³⁷ and Prescott et al.,³⁸ respectively. ChIP-seq data from craniofacial data in Wilderman et al.³⁹ comprise processed formats, including imputed signals, peaks, and segmentation data. In order to ensure comparability among the three data sources, computational processing of ChIP-seq data as published in Wilderman et al.³⁹ (QC, alignment, peak calling, epigenetic imputation, chromatin segmentation) was adopted to the hNCC/cNCC bioinformatics pipeline, as described in the Supplemental information and Table S5.

Chromatin imputation and segmentation

To obtain uniform datasets, chromatin imputation followed by chromatin state segmentation was performed. First, H3K9me3 and H3K36me3 marks in hNCCs/cNCCs were imputed using ChromImpute (v.1.0.1),⁴⁸ based on 127 cell types from the Roadmap Epigenome Project.³⁴

Imputed hNCC/cNCC signal files for each individual chromosome and each chromatin mark were binarized, and segmentation was performed using the core+K27ac 18-state chromatin model provided by Roadmap with ChromHMM⁴⁹ to predict 18 chromatin states. Because of the low number of chromatin marks measured in the NCC samples, epigenetic imputation issues, and the higher risk of batch effect between hNCCs, cNCCs, and CT, we adopted a robust strategy and condensed the 18 generated states into eight states, based on Roadmap definition: three active states (transcription starting sites [TSS], transcribed sites, and enhancers [Enh]), one bivalent state (Poised Enh/bivalent TSS), three repressed states (Heterochromatin, Repressed PolyComb sites, Zinc finger genes/Repeats), and one quiescent state (Quies). Potential batch effects were analyzed using principal-component analysis (PCA) and hierarchical clustering of Pearson correlation coefficients.

Other datasets

To identify genome-wide regulatory genomic units, we used TADs from human embryonic stem cells (hESCs) (H1 cell line) as provided by the Ren Lab.⁴⁰ Protein-coding genes were extracted from UCSC genome browser (hg19) and were mapped to TADs using positional information. TADs containing an nsCL/P risk locus were defined as GWAS_{TAD} region. Based on previous evidence for complex regulatory interactions within one TAD, we considered all genes from the GWAS_{TAD} region as potential candidate genes for downstream effects of the associated variants in the $r^2 \geq 0.6$ region. Expression data from NCCs (two replicates of day11hNCC [GEO: GSE121428] and three replicates of passage2hNCC [GEO: GSE108521]) were retrieved from Lausch et al. (GEO: GSE108522).⁵⁰ For the comparison of genes in TADs of nsCL/P risk loci and genes expressed in NCCs, we used the average RNA-seq Fragments Per Kilobase Million (FPKM) across five samples. To identify functional links between different regulatory features (e.g., DNA-DNA interactions of enhancers and TSS) at specific risk loci, we accessed pChIP-C *cis*-interaction data collected in hESCs (GEO: GSE86821).⁴¹

Translation of genetic associations into tissue- and time-point-specific regulatory effects at a systematic level

Enrichment analyses using GREGOR

Based on chromatin segments obtained from hNCCs, cNCCs, and CT, we used GREGOR (Genomic Regulatory Elements and GWAS Overlap Algorithm)⁵¹ to evaluate the enrichment of significant

SNPs from the MAiC data in the available regulatory features (i.e., eight predicted chromatin states). As described in the Supplemental information, a set of samples from the Roadmap Epigenomics project (comprising both fetal and adult tissue samples) was selected as an independent dataset for comparison. As input, we used MAiC nsCL/P variants with $p \leq 0.001$ without additional variants in LD ($n = 22,999$); this threshold was selected to balance between adequate statistical power and true-positive association signals.

CT- and NCC-specific active chromatin sites

To examine specific effects in either NCCs or CT, we filtered in the chromatin segmentation datasets for active chromatin sites (TSS, Enhancer or transcribed sites) in NCCs that are repressed/quiet (Quiescent, Biv_TSS_pois_enh, ReprPC, Heterochromatin) in CT and vice versa. For robust observations, we only trust in a chromatin state if it is present in both NCC samples (hNCCs, cNCCs) or in five of the six CT (CS13, CS14, CS15, CS17, CS20, 10wpc) samples. To account for biases in length associated with batch effects, active sites were only retained if they had a distance of ≥ 500 bp to any chromatin segment of opposite activity status in the other cell system/tissue. In the following, we combined the specific active chromatin sites with MAiC associations and TAD data to filter for TADs with high density of strong associated genetic variants ($p_{\text{MAiC}} \leq 5 \times 10^{-5}$) in specific active chromatin sites at new and known nsCL/P GWAS risk loci.

Characterization of nsCL/P risk variants and candidate gene prioritization in context of epigenetic mid-facial timeline

For comprehensive insights in regulatory mechanisms at nsCL/P risk loci, we finally integrated all available genetic and functional data (MAiC associations, GWAS_{TAD}- and $r^2 \geq 0.6$ -region boundaries, NCC- and CT-specific active chromatin sites, chromatin segmentation tracks, and pChIP-C *cis* interactions). Based on this approach, we attempt to prioritize genetic variants with regulatory effect and potential downstream target genes and to detect relevant regulatory elements specific for the early (hNCC/cNCC) or later mid-facial development (CT).

Results

MAiC identifies five novel risk loci

The MAiC dataset was generated by combining GWAS data from three previous studies (Bonn,¹⁸ GENEVA,²⁴ POFC¹⁷), following the exclusion of overlapping individuals and extensive QC. The final dataset comprised 1,247 nsCL/P cases, 2,879 controls, and 2,699 case-parent trios of multiple ethnicities, and ~ 7.74 million SNPs. The p value distribution was consistent with a multifactorial inheritance (Figure 1B; $\lambda = 1.07$). A set of 1,375 SNPs achieved genome-wide significance ($p < 5 \times 10^{-8}$; Figure 1C). Analysis of established nsCL/P risk loci in MAiC revealed genome-wide significant SNPs at 25 of the 40 regions. These 25 regions comprised 22/26 loci that were previously identified in GWASs based on largely European samples and 3/14 loci reported in individuals from the Chinese population.^{13,22} At all other nsCL/P risk loci ($n = 15$), nominal significance ($p < 0.05$) was observed for individual

variants that were in strong LD ($D' > 0.8$) with the respective lead SNP (Table S2).

Importantly, the MAiC analyses also identified five novel risk loci ($p < 5 \times 10^{-8}$), thus increasing the number of identified nsCL/P GWAS risk loci to 45. These novel loci were located at chromosomes 1p36.13 (sentinel variant rs34746930), 5p12_{FGF10} (rs60107710), 5q13.1_{PIK3R1} (rs6449957), 7p21.1 (rs62453366), and 20q13.12 (rs3091552; Table 1). Consistent with previous findings on risk variants for nsCL/P and other complex traits,²⁹ these lead variants map to non-coding regions that are adjacent to candidate genes with functions during facial development, such as *CAPZB*⁵² and *NBL1*⁵³ (both at 1p36) and *EYA2*⁵⁴ (at 20q13; Supplemental text; Figures S2–S6). To identify population-specific effects, a sub-analysis was performed in individuals from Central Europe (MAiC_{Euro}; $n = 562$ cases, 2,051 controls, and 955 case-parent trios). No additional risk loci were identified at the level of genome-wide significance (Figure S7; Table S2). Using this European case-control cohort and LD score regression,⁴⁵ SNP-based heritability was estimated as $h^2 = 28\% \pm 0.1\%$. This confirmed previous heritability estimates obtained using the Bonn cohort only.²¹

Gene-based analyses suggest nsCL/P candidate genes outside of GWAS risk loci

Using MAiC summary statistics and MAGMA,⁴⁶ gene-based analyses yielded 1,357 genes with nominal significance ($p < 0.05$; Figure S8A). A total of 25 genes reached test-wide significance ($p < 2.79 \times 10^{-6}$; Table S6). Of these, 23 map to known GWAS risk loci. For some of these 23 genes, functional evidence strongly supports their involvement in nsCL/P (e.g., *IRF6*,⁵⁵ *TP63*⁵⁶). This analysis also suggested novel candidate genes at GWAS risk loci, such as *ARID3B*. In mice, the gene *Arid3b* is expressed in cranial mesenchyme structures and has been shown to interact with *Mycn*, which is encoded by a strong candidate gene at another nsCL/P risk locus.^{57,58} Two genes with a significant burden of common variants mapped outside all known GWAS risk loci. These genes, *BTN3A3* ($p_{\text{gene}} = 6.96 \times 10^{-7}$) and *BTN3A1* ($p_{\text{gene}} = 2.44 \times 10^{-6}$; Figure S9A), are both located at chromosome 6p22.2, and previous research found that *BTN3A3* showed differential expression in the lip tissue of CL/P phenotypic subgroups.⁵⁹ In MAiC_{Euro}, the gene-based analysis revealed 11 genes with test-wide significance (Figure S8B; Table S7), including three novel candidate genes (*LIMCH1*, *MSX2*, and *STRA13*; Figures S9B–S9D). Overall, 41 genes yielded $p < 10^{-5}$ in one of the two analyses.

We also analyzed a set of 13 previously identified nsCL/P candidate genes with: (1) a significant enrichment of low-frequency variants (four genes),⁶⁰ (2) an autosomal-dominant inheritance pattern in multigenerational families (four genes),⁶¹ or (3) an enrichment of rare coding variants (five genes).⁶² Of these, 12 genes were present in the analysis set. Two of these 12 genes approached test-wide significance: *PRTG* ($p = 8.44 \times 10^{-5}$) and *CTNND1* ($p = 2.17 \times$

10^{-5} ; Table S8). These observations indicate that in at least a subset of genes, both common and rare variations, contribute to nsCL/P.

Genes located in TAD regions of nsCL/P GWAS loci are enriched in developmental pathways

Accumulating evidence suggests that most regulatory interactions occur within TAD modules.^{63,64} Therefore, genes located within TADs represent candidates for the downstream effects of the associated common risk variants. To identify molecular processes of relevance to nsCL/P, for each of the 45 risk loci, GWAS_{TAD} regions were defined, based on the extent of the respective TAD in hESC data.⁴⁰ In total, 407 genes were identified within the respective TADs (GWAS_{TAD} genes, range 1 to 29 genes per locus; Table S9). Enrichment analysis using MAGMA yielded test-wide significant ($p_{\text{adj}} \leq 0.05$) results for 287 GO terms (Table S10). The most significant enrichments were observed for “tissue development” ($p_{\text{adj}} = 8.34 \times 10^{-9}$), “epithelium development” ($p_{\text{adj}} = 8.82 \times 10^{-9}$); and “appendage development” ($p_{\text{adj}} = 7.92 \times 10^{-8}$; Figure S10). Together with additional significant terms, such as “embryo development,” “tube development,” and “ear development,” these observations suggest the existence of common pathways for nsCL/P and other processes of organogenesis during embryonic development.

We then prioritized genes expressed in NCCs by adding available RNA sequencing (RNA-seq) data from hNCCs.⁶⁵ In total, 240 of the 407 GWAS_{TAD} genes were expressed in NCCs, with strong expression being observed for a subset of 12 genes (≥ 200 fragments per kilobase mapped; Table S9). Of these, at least two have been previously implicated in NCC migration processes (*CAPZB*,⁵² *TPM1*⁶⁶). These 240 NCC-expressed genes showed a substantial overlap in significant GO terms compared with the analysis of all 407 GWAS_{TAD} genes (233 out of 287 pathways; Figure 2A; Table S11). Of those 233 pathways, 157 pathways showed stronger enrichment in the subset of NCC-expressed GWAS_{TAD} genes, the strongest of which represent cellular processes (Figure S10; Table S12). Among pathways that were exclusive to GWAS_{TAD} genes expressed in NCCs ($n = 106$), both regulatory processes and metabolic pathways were enriched. In contrast, pathways specific to GWAS_{TAD} genes that were not expressed in NCCs ($n = 54$) included “keratinocyte proliferation” and “epidermis development,” a finding that is consistent with the substantial contribution of the epithelial lineage to nsCL/P.⁵⁶

We next addressed the potential etiological overlap between nsCL/P and other common phenotypes that might contribute to the adverse health outcomes observed in nsCL/P. We retrieved association signals for each of the 45 lead SNPs in MAiC from large-scale genetic studies, using the GWAS ATLAS.⁴⁷ At the time of analysis (November 22, 2019), this resource comprised 4,756 GWASs on 3,302 unique traits. While all of the 45 variants were available in the atlas, only 19 showed at least one significant SNP-trait association when corrected for the number of GWASs and

Table 1. Novel risk loci for nsCL/P identified in MAiC

Locus	Lead variant	Position ^a	Allele 1/allele 2 ^b	p value	RR ^c	95% CI
1p36.13	rs34746930	19,781,724	<u>C</u> /G	4.19×10^{-8}	1.30	1.18–1.43
5p12 _{FGF10}	rs60107710	44,577,755	<u>A</u> /G	3.50×10^{-8}	1.39	1.24–1.57
5q13 _{PIK3R1}	rs6449957	67,483,732	<u>T</u> /C	6.59×10^{-9}	1.21	1.13–1.29
7p21.1	rs62453366	20,747,107	G/ <u>T</u>	7.83×10^{-9}	0.77	0.70–0.84
20q13.12	rs3091552	45,440,006	<u>C</u> /G	1.31×10^{-9}	1.38	1.22–1.47

nsCL/P, non-syndromic cleft lip with or without cleft palate; MAiC, meta-analysis in clefting; RR, relative risk; CI, confidence interval. Gene names in subscript distinguish novel associated regions from independent risk loci at the same chromosomal band.

^aPosition according to hg19.

^bRisk allele is underlined.

^cRR provided for allele 1.

loci ($p < 2.33 \times 10^{-7}$; overall number: $n = 219$; Table S13). These associations reflect 35 collapsed traits across 12 domains, including height, bone mineral density, hair color, and body mass index (Table S14). Eighteen traits showed associations with at least two distinct nsCL/P risk loci. Interestingly, for some traits, the direction of effect differed between individual loci (e.g., height and bone mineral density), while for other traits, the direction of effect was consistent (e.g., hypothyroidism, glomerular filtration rate, and hair color; Figure 2B).

NsCL/P-associated variants are enriched in multiple chromatin states of mid-facial development

Recent analyses in human embryonic CT³⁹ demonstrated both a significant enrichment of lead SNPs from earlier nsCL/P GWAS in active enhancers and the presence of mid-facial specific regulatory elements. To extend this work, we incorporated data from two NCC states in order to generate a unified mid-facial development resource of chromatin modifications (Figure S1). We retrieved data on ChIP-seq from hNCCs³⁷ and cNCCs³⁸ and applied the data analysis pipeline used by previous authors for computational analyses of ChIP-seq data from CT.³⁹ We observed strong inter-sample correlations between chromatin mark and developmental stage (Figures S11 and S12). The integration of 127 non-facial samples from Roadmap³⁴ revealed local clustering of NCCs and CT along a hierarchical axis comprising hESCs, induced pluripotent stem cells (iPSCs), and iPSC-derived cells (Figure S13). Here, the most tissue-specific pattern was observed for H3K27ac (Figure S14). Similar to a previous finding for CT,³⁹ non-facial fetal tissue samples (such as brain, kidney, and lung) clustered distinctly from NCCs (Figure S14), thus emphasizing the limited utility of many public resources for the interpretation of genetic findings in facial disorders.

Next, we generated robust chromatin segments in NCCs using ChromHMM.⁶⁷ Together with segmentation data from CT and Roadmap, chromatin segments were condensed to eight categories in order to increase the robustness of the subsequent analyses (Figure S15; Table S15). We then analyzed the positional overlap of all variants with $p_{\text{MAiC}} < 0.001$ in the eight chromatin states across NCCs and CT (SNP_{0.001_nsCL/P} $n = 22,999$), and

compared this to a matched set of non-associated SNPs (SNP_{control_nsCL/P} $p > 0.1$). The results showed that 23% of the nsCL/P variants (SNP_{0.001_nsCL/P}) mapped to active chromatin states, while 14% mapped to either bivalent or repressed chromatin states (Figure 3A). This enrichment was significantly higher compared to the control SNPs, where 16% and 11% of variants mapped to active, or to bivalent/repressed, chromatin states, respectively ($p < 10^{-16}$, Fisher's exact test).

To delineate associations of specific chromatin states along the time series, enrichment was tested using GREGOR.⁵¹ For each of the two SNP sets, every hNCC/cNCC/CT sample was tested, together with 11 randomly selected Roadmap samples (both fetal and adult). A significant enrichment for SNP_{0.01_nsCL/P} was observed in most of the samples/chromatin states (Figure 3B; Table S16), as compared to SNP_{control_nsCL/P} (Figure S16; Table S17). While the fold enrichment (FE) was similar for NCCs and CT in six of the eight chromatin states (such as those related to active transcription; Figures 4A–4D; Figure S17), considerable differences in enrichment between NCC and CT samples were observed in chromatin states “active enhancers” and “poised enhancers/bivalent TSS.” In both states, NCCs displayed a stronger enrichment than CT samples. For enhancers, the mean FE (FE_{Mean}) in NCCs was 1.64 ($p_{\text{Mean}} = 4.36 \times 10^{-86}$, average of p_{GREGOR}), compared with FE_{Mean} = 1.43 in CT ($p_{\text{Mean}} = 8.09 \times 10^{-22}$). For “poised enhancers/bivalent TSS,” the corresponding values were FE_{Mean} = 1.65, $p_{\text{Mean}} = 3.39 \times 10^{-20}$ in NCCs, compared with FE_{Mean} = 1.39, $p_{\text{Mean}} = 4.74 \times 10^{-4}$ in CT. These results may have been driven in part by the heterogeneous composition of the CT samples. However, the specific enrichment pattern observed in two out of eight chromatin states suggests a distinct biological underpinning. Overall, these data confirmed previous findings of an overrepresentation of nsCL/P lead variants in enhancer marks^{21,39} and extended this enrichment toward additional common variants and annotations.

A subset of nsCL/P-associated SNPs show distinct regulatory effects

To extend the investigation of the contribution of regions with differing regulatory profiles in NCCs and CT, we

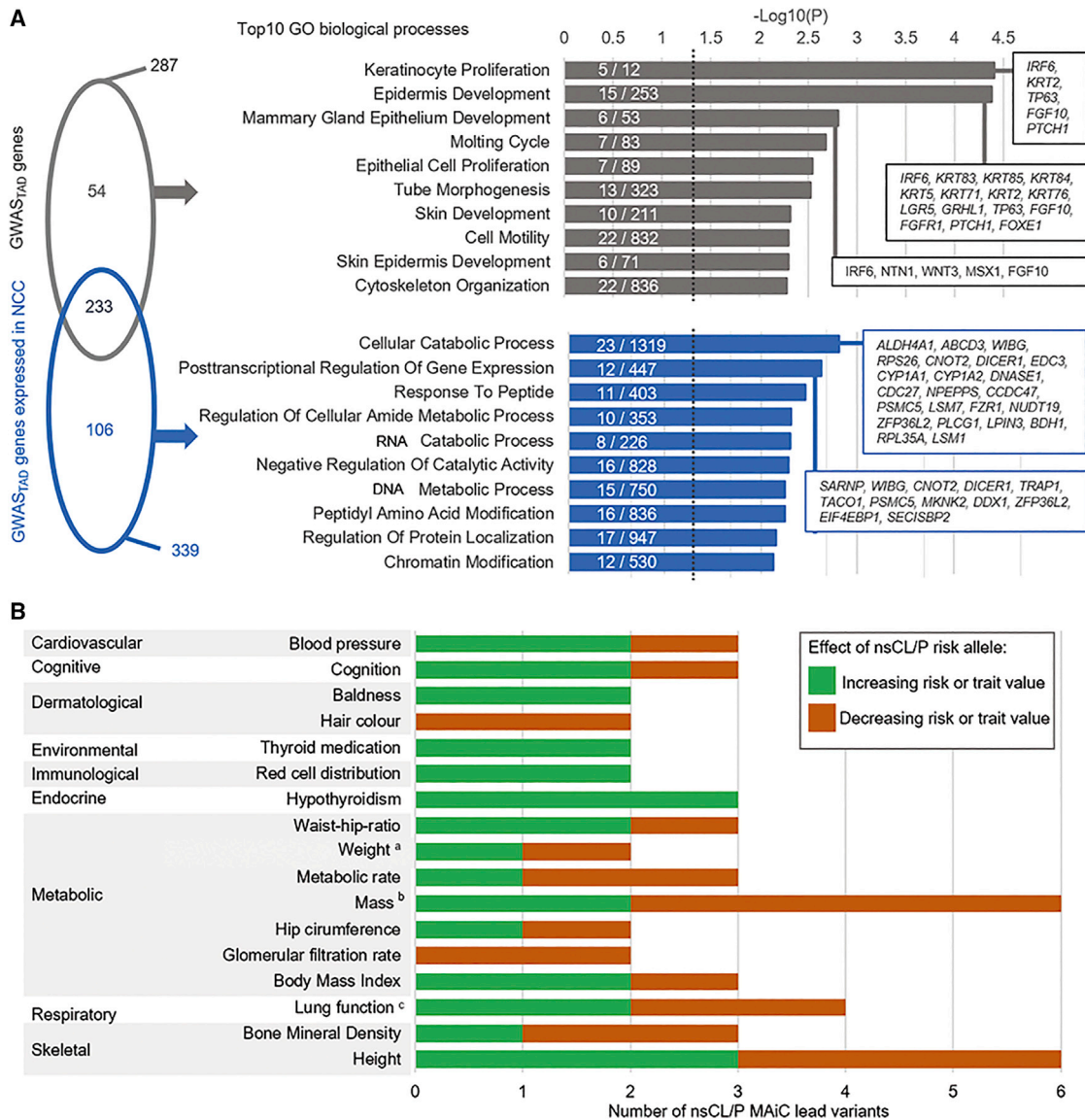


Figure 2. Systematic assessment of 45 risk loci for nsCL/P

(A) Enrichment analyses of biological processes. Enrichment of genes located at risk loci identified by genome-wide association studies (GWAS_{TAD} genes, $n = 407$, gray) and the subset of genes expressed in neural crest cells ($n = 240$, blue) were calculated using MAGMA. Left panel: While most of the associated pathways overlapped both datasets, a subset of terms was distinctly enriched in one of the groups. Right panel: Bars represent the top 10 of each specific enrichment ($p_{adj} \leq 0.05$). Numbers reflect nsCL/P risk genes/total number of genes in the respective gene ontology (GO) term. For the most strongly associated pathways, gene names are provided in the respective box. (B) Pleiotropic effects of lead variants. For the lead variant of each of the 45 nsCL/P risk loci, associations with common traits were retrieved from the GWAS ATLAS. Associations with at least two risk loci were observed for 17 traits from 12 domains (y axis). Bar colors represent direction of effects. ^aIncluding birth weight. ^bIncluding multiple mass-related measurements. ^cLung function as measured by Forced expiratory volume (FEV)₁ or FEV₁/Forced vital capacity (FVC) ratio.

created genome-wide maps of active chromatin sites for both NCCs and CT. A total of 9,897 regions (encompassing 26.67 Mb) with active chromatin states in NCCs (TSS, enhancer or transcribed sites) were inactive in CT (quiescent, repressed, or bivalent; termed NCC-specific active sites). Similarly, 6,189 regions (29.37 Mb) were active in CT but inactive in NCCs (CT-specific active sites). The integration of MAiC association data revealed 62,084 genetic variants that map in NCC-specific active sites. Of these, 4,022 had $p_{MAiC} \leq 0.05$. Similarly, 72,556 variants (4,834 of which had $p_{MAiC} \leq 0.05$) mapped to CT-specific active

sites. In each of the groups of NCC-specific and CT-specific active sites, the p value distribution differed significantly from that expected, with a significant enrichment of association signals being observed at the lower tail of the distribution (Figure 5A).

Filtering for the subset of SNPs with $p_{MAiC} \leq 5 \times 10^{-5}$ identified 112 SNPs that mapped to either NCC-specific (51 variants), or CT-specific active regions (61 variants; Table S18). These were distributed over 39 TADs, which encompassed both known nsCL/P risk loci ($n = 19$; e.g., chromosomes 1p22 [Figure S18] and 2p24.2 [Figure S19]) and

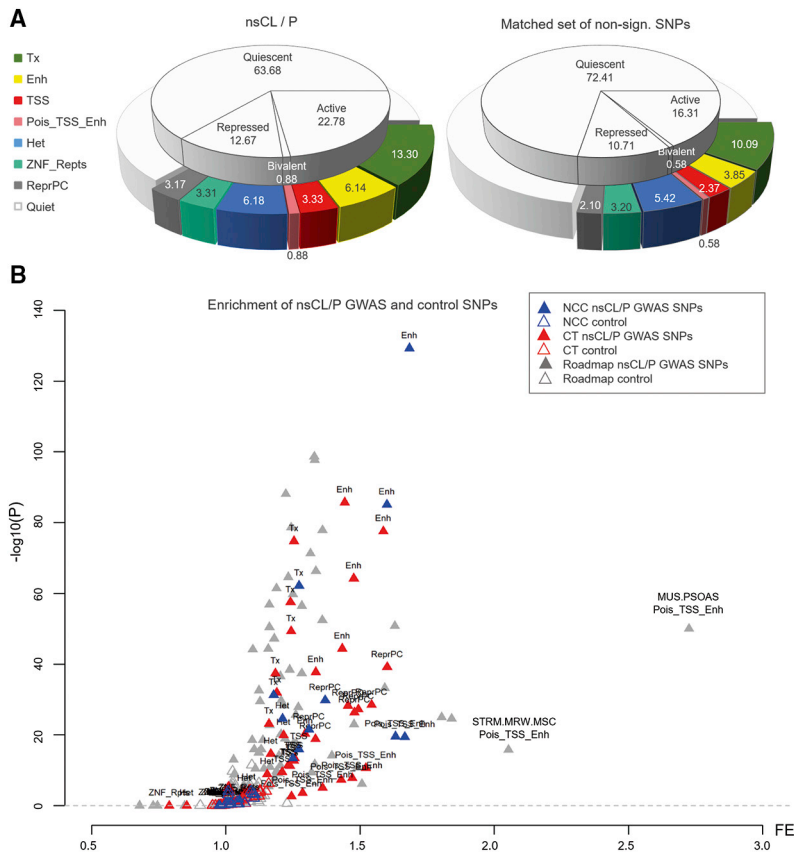


Figure 3. Association of MAiC across epigenetic annotations

For all enrichment analyses, two sets of single-nucleotide polymorphisms (SNPs) were designed: (1) set of MAiC risk variants, at $p_{\text{MAiC}} \leq 0.001$ ($n = 22,999$), and (2) a size-matched control set, comprising non-associated SNPs ($p_{\text{MAiC}} > 0.1$) with similar allele-frequency distribution.

(A) Overall enrichment analysis. For each group, the fraction of SNPs represented in different chromatin annotations of mid-facial development was assessed, without discriminating between NCCs and craniofacial tissue (CT).

(B) Overview of enrichment in NCCs and CT. Enrichment of nsCL/P risk variants in eight chromatin states for each sample (hNCCs, cNCCs, and CT, plus a set of 11 Roadmap samples). p values were calculated using GREGOR.⁵¹ Abbreviations: TSS, transcription starting site; Enh, enhancer; ReprPC, repressed PolyComb; Tx, transcribed sites; Het, Heterochromatin; TxFlnk, transcribed sites at gene 5' and 3'; Pois_TSS_Enh, poised enhancers and bivalent TSS; ZNF_Rpts, Zinc finger genes and repeats; FE, fold enrichment. Abbreviations of tissues as provided by Roadmap.³⁴

regions with suggestive evidence for association ($n = 20$; e.g., chromosome 4p13 [Figure S20]). Interestingly, at six loci (e.g., chromosomes 1q32.1 [Figure S21] and 15q24.1 [Figure S22]), at least two associated variants in LD were located in different specific elements (Table S19). This represents a significantly higher enrichment than expected and suggests that individual variants of risk haplotypes might affect the regulatory architecture at different stages of craniofacial development (Figure 5B).

Finally, we assessed how novel hypotheses on nsCL/P pathogenesis can be generated from the systematic integration of data concerning: (1) statistical associations (MAiC), (2) chromatin modifications over time (mid-facial time-series), and (3) pChI-C *cis*-interactions.⁴¹ Examples from two loci are described here. First, at 5q13_{PIK3R1}, the lead variant (rs6449957, $p_{\text{MAiC}} = 6.59 \times 10^{-10}$) is located within an active region upstream of *PIK3R1*. This region shows evidence of being transcribed but lacks any RefSeq annotation, which might point toward a transcribed enhancer or an as-yet-undetected transcript. PChI-C data indicate *cis* interactions with *PIK3R1* and *MAST4*, both of which are expressed in hNCCs. In addition, another variant in strong LD (rs921792, $p_{\text{MAiC}} = 1.17 \times 10^{-5}$) maps to a putative enhancer that is detected in both NCCs and CT (Figure 5C). As a second example, at 13q32.2 (lead variant rs2763950, $p_{\text{MAiC}} = 3.03 \times 10^{-6}$, intronic in *CLYBL*), interactions were observed between the region around the lead variant and the genes *ZIC2*, *ZIC5*, and *GGACT*. While

some variants (including rs2763934 with $p_{\text{MAiC}} = 6.53 \times 10^{-7}$) map to a craniofacial active element near the *CLYBL* gene promoter, additional variants (including rs4525350 with $p_{\text{MAiC}} = 6.39 \times 10^{-6}$) map to several more distantly located NCC-specific enhancers.

Based on pChI-C data, our data indicate that in NCCs, risk variants might affect *ZIC2* and *ZIC5* expression. This hypothesis is further supported by the finding of active transcription sites in NCCs and a bivalent state in embryonic and adult tissues. A plausible hypothesis is that, at later time points of development, additional variants mapping to other enhancer elements act on *GGACT*, as suggested by the presence of transcribed sites in CT. Notably, the transcript region of *CLYBL* itself has limited evidence for active transcription across all analyzed stages of mid-facial development, despite the presence of some active marks in the promoter region (Figure 5D).

At other loci, our data provide evidence for the presence of tissue-specific gene isoforms (e.g., 4p13-locus; Figure S20), or a second, novel candidate gene at previously reported loci. For example, at chromosome 1p22, our data suggest that the previously identified gene *ARHGAP29*³⁶ is a target gene with CT-specific expression and highlight *ABCD3* as novel candidate gene (Figure S18). The data also suggest complex promoter-promoter interactions involving all genes at this locus (*ARHGAP29*, *ABCD3*, and *ABCA4*). Interestingly, the MAiC top-associated variant at 1p22 (rs35298667, $p_{\text{MAiC}} = 6.86 \times 10^{-16}$) has putative enhancer function and maps to the “E2” element, whose functional role in nsCL/P was confirmed in previous research.³² At another locus (1q32.1), we found that *SERTAD4* is a CT-specific target gene, while the established causal gene *IRF6* was marked as bivalent, which

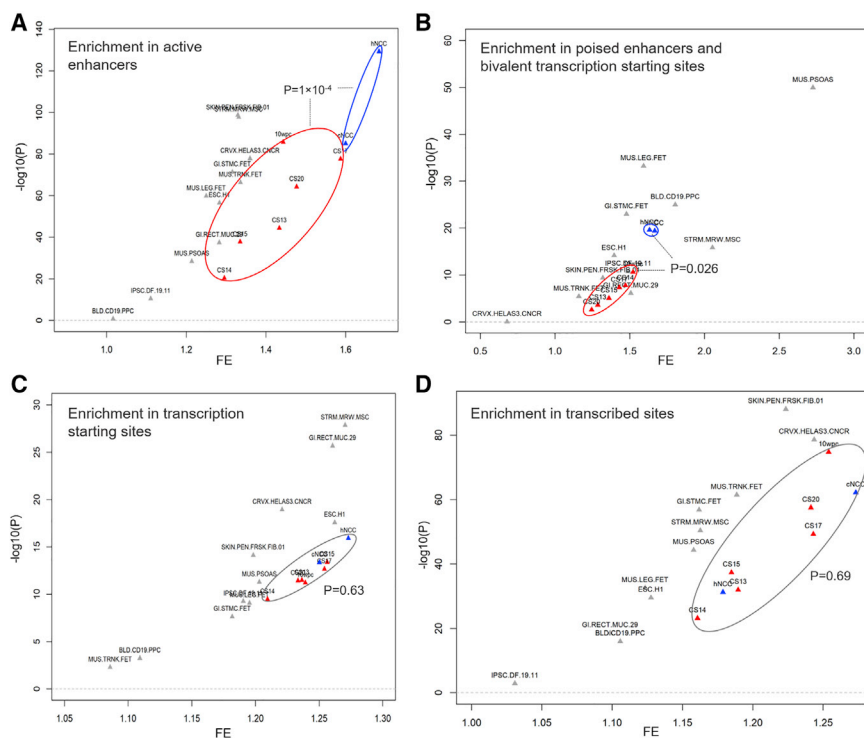


Figure 4. Association of MAiC variants in distinct chromatin states

(A–D). Individual enrichment results for MAiC risk variants in four chromatin states. p values represent the difference in enrichment between NCCs and CT. Abbreviations: TSS, transcription starting site; Enh, enhancer; ReprPC, repressed Poly-Comb; Tx, transcribed sites; Het, Heterochromatin; TxFlnk, transcribed sites at gene 5' and 3'; Pois_TSS_Enh, poised enhancers and bivalent TSS; ZNF_Rpts, Zinc finger genes and repeats; FE, fold enrichment. Abbreviations of tissues as provided by Roadmap.³⁴

is consistent with its established activity in epithelial tissue^{55,68} (Figure S21). Taken together, these results will inform future functional studies into nsCL/P and underscore the importance of the thorough genomic annotation of relevant cells and tissues.

Discussion

Here, we report on a data-driven approach that generated novel insights into the etiology of nsCL/P. At the genetic level, we identified five novel risk loci via the large-scale meta-analysis of common genetic variation. This large genome-wide resource empowered systematic analyses at the gene and pathway levels and implicated novel molecular players in nsCL/P. Our analysis of pleiotropic effects on other common traits revealed a substantial positional overlap with traits such as height and bone mineral density. At some loci, associated variants showed opposite directions of effect, which indicates their contribution to distinct pathways. We have provided examples of how this resource is useful in terms of translating statistical associations into biological insights and illustrated its potential for further analyses of facial disorders and traits.

While our results are based on a multiethnic cohort, this still comprises a substantial contribution from the European population. Still, we captured associations at all loci that had been previously reported in distant ethnicities, such as the Chinese population.¹³ Although these observations suggest that nsCL/P might show less locus heterogeneity than is the case for other common diseases, allelic heterogeneity is likely to contribute in part to the lack of replication

observed at some loci in previous studies. Also, the integration of genetic and chromatin segmentation data might have been biased by the European background of both the genetic and epigenetic maps. Despite some initial evidence that methylation patterns show population-specific components,^{69,70} few studies to date have performed systematic analyses of how maps of chromatin accessibility (in particular in mid-facial development) vary across populations. Future studies are required to determine whether population-specific risk variants from non-European populations show differing enrichment patterns from those observed in the present study and to identify additional pleiotropic effects that are present at other risk haplotypes in other populations. Importantly, to address these issues, future meta-analysis should also include recent GWAS data (e.g., from Sub-Saharan Africans⁷¹ and Colombians⁷²). In addition, our analyses were performed for nsCL/P as the central trait. Previous studies have generated evidence of an (albeit incomplete) etiological overlap between the various nsCL/P subtypes (e.g., cleft lip, and cleft lip with cleft palate) and the genetic heterogeneity of other types of orofacial clefting (e.g., cleft palate only).^{21,23,73,74} Application of our integrative approach to the investigation of cleft subtypes will facilitate understanding of their individual etiologies, an issue that was beyond the scope of the present study.

One major feature of our approach was that it combined previous individual data into one joint map of epigenetic chromatin segments of NCCs and CT. This will be highly useful in terms of the future interpretation of associations in facial disorders/traits. However, due to limited availability of datasets from other cell types, such as human embryonic epithelium, this map does not comprehensively capture all biological contributors to human craniofacial development. Furthermore, our joint analysis of the different CT stages may have overlooked some effects within single stages of CT. Nonetheless, the data obtained at individual loci add to increasing evidence that for nsCL/P development, risk loci have a complex regulatory

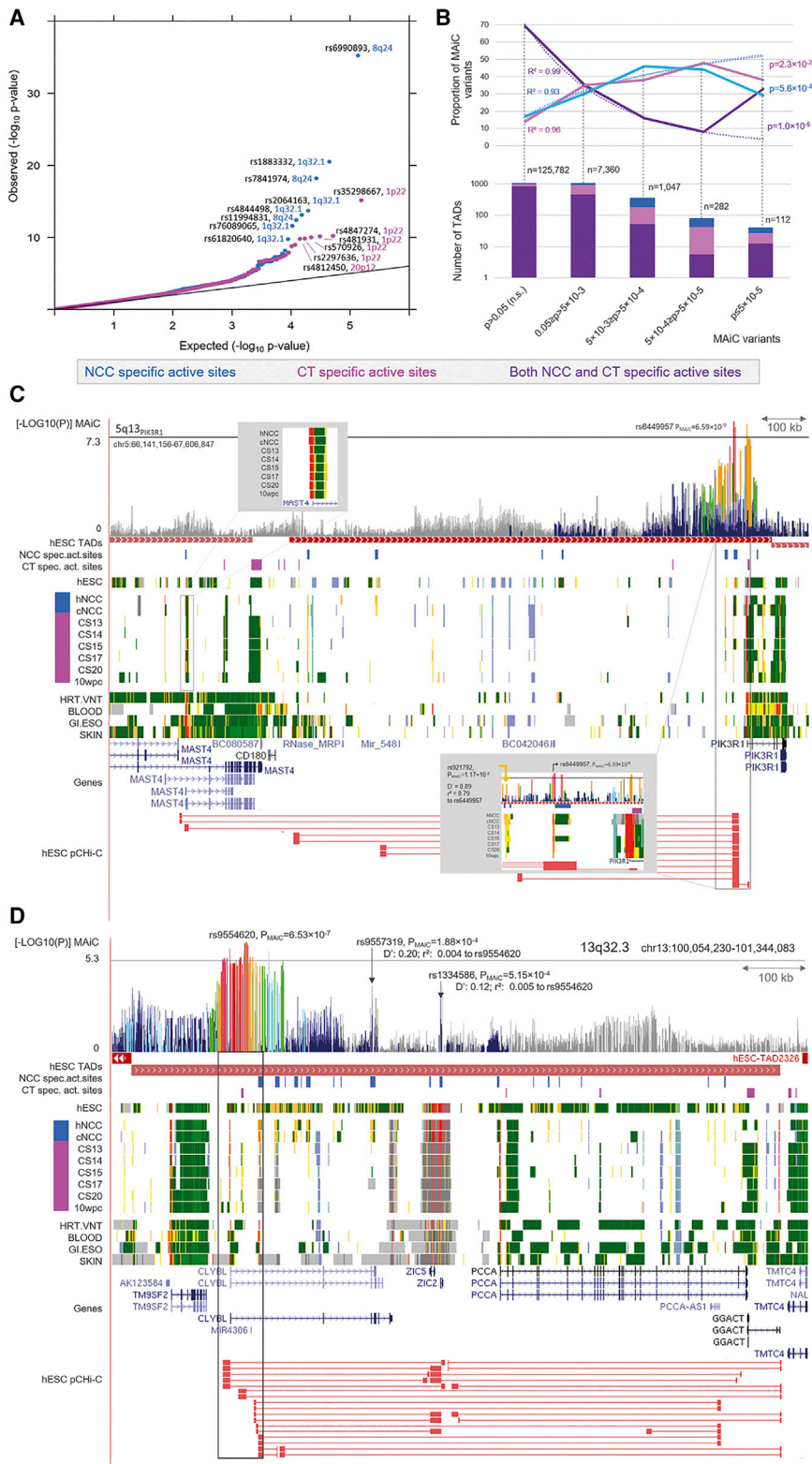


Figure 5. Interpretation of MAiC association results

(A) Quantile-quantile plot of specific active sites. P_{MAiC} values (as $-\log_{10}$) of SNPs located in NCC-specific ($n = 62,084$; blue) or CT-specific ($n = 72,556$; pink) active sites are plotted against expected p values. In both datasets, a significant enrichment of associated risk variants was observed.

(B) Distribution of risk variants in specific active sites. Variants located within NCC- and CT-specific regions were retrieved at different P_{MAiC} cutoffs and aggregated per topologically associated domain (TAD, numbers in lower panel). TADs were classified according to whether the variants map uniquely to NCC-active elements (blue), CT-specific elements (pink), or both (purple). The distribution largely followed the expected logarithmic distribution. However, for a substantial number of loci, different associated SNPs (at $p < 5 \times 10^{-5}$) mapped to both NCC- and CT-specific sites within one TAD.

(C and D) Regulatory architecture at selected loci. Based on the extent of the TAD around the respective lead variant and variants in LD ≥ 0.6 (shown in gray framed box), different layers of data were aggregated and are represented for risk loci 5q13_{PIK3R1} (C) and 13q32.3 (D). Tracks include (top-down): MAiC p values with color code based on LD to respective top variants; extent of NCC-specific (blue) and CT-specific (pink) sites; chromatin segmentation data from hNCCs, cNCCs, CT (color code as in Figure 3), and selected samples from Roadmap; RefSeq gene positions; and promoter capture (pC) Hi-C *cis*-interactions collected in hESCs.

ARHGAP29^{76,77} and *IRF6*²⁵), or experimental evidence (e.g., *PAX7*⁷⁸). While we here focused on an *in silico* approach, we hope that the results will empower further experimental investigations of specific risk variants that were highlighted among the set of associated variants. Using the joint pipeline, we will continue to update our resource as chromatin marks become available from additional human tissues and/or cell systems of relevance to mid-facial development. In addition, the map will be refined through the use of single-cell technologies

architecture, and several genes at single loci might be of relevance across the different time points of craniofacial development. Notably, several of the genes prioritized by our systematic approach have obtained independent support by other studies, for instance clefting syndromes (e.g., *TP63*, EEC syndrome⁷⁵), resequencing studies (e.g.,

ologies in order to resolve the issue of tissue heterogeneity encountered in the present study. Finally, the integration of other layers of genetic information, such as rare variants identified by whole-exome or -genome sequencing in cleft cohorts, will further increase our understanding of the etiology of craniofacial development and disease.^{61,79}

Data and code availability

Original data for genetic and functional analyses in the paper is available as follows: dbGaP (dbGaP: phs000094 and phs000774), GEO (GEO: GSE28874, GSE70751, and GSE97752), and Zenodo (DOI 10.5281/zenodo.3724148). The NCC- and CT-specific active sites generated during this study are available at Zenodo (DOI 10.5281/zenodo.3911187).

Supplemental information

Supplemental information can be found online at <https://doi.org/10.1016/j.xhgg.2021.100038>.

Acknowledgments

We thank Markus M. Nöthen and Andreas Bunes for helpful discussions on the manuscript and Carlo Maj for data management. This work was supported by the German Research Foundation (DFG; LU 1944/3-1, to K.U.L.).

Declaration of interests

The authors declare no competing interests.

Received: January 19, 2021

Accepted: May 27, 2021

Web resources

ANNOVAR, <https://wglab.org/software/9-annovar>
Bowtie2, <http://bowtie-bio.sourceforge.net/bowtie2/index.shtml>
ChromHMM, <http://compbio.mit.edu/ChromHMM/>
ChromImpute, <http://www.biolchem.ucla.edu/labs/ernst/ChromImpute/>
core 15-state chromatin model, <https://egg2.wustl.edu/roadmap/data/byFileType/chromhmmSegmentations/ChmmModels/>
FastQC, <https://www.bioinformatics.babraham.ac.uk/projects/fastqc>
FUMA, <https://fuma.ctglab.nl/>
GREGOR, <http://csg.sph.umich.edu/GREGOR/>
GTEx, <https://www.gtexportal.org/home/>
GWAS Atlas, <https://atlas.ctglab.nl/>
Hi-C data from Bing Ren Lab, <http://chromosome.sdsc.edu/mouse/hi-c/download.html>
IMPUTE2, http://mathgen.stats.ox.ac.uk/impute/impute_v2.html
KING: Kinship-based Inference for GWAS, <https://www.kingrelatedness.com/>
LDlink, <https://ldlink.nci.nih.gov/?tab=home>
LDSR, <http://ldsc.broadinstitute.org/ldhub/>
MACS2, <https://github.com/macs3-project/MACS>
MAGMA, <https://ctg.cncr.nl/software/magma>
METAL, <http://csg.sph.umich.edu/abecasis/metal/>
MSigDB, <https://www.gsea-msigdb.org/gsea/msigdb/index.jsp>
PhantomPeakQualTools, <https://github.com/kundajelab/phantompeakqualtools>
Roadmap, <https://egg2.wustl.edu/roadmap/data/byFileType/signal/consolidated/macs2signal/pval/>
UCSC, <http://genome.ucsc.edu/>

References

1. Claussnitzer, M., Dankel, S.N., Kim, K.-H., Quon, G., Meuleman, W., Haugen, C., Glunk, V., Sousa, I.S., Beaudry, J.L., Puvion-Rodan, V., et al. (2015). FTO Obesity Variant Circuitry and Adipocyte Browning in Humans. *N. Engl. J. Med.* 373, 895–907.
2. Xue, A., Wu, Y., Zhu, Z., Zhang, F., Kemper, K.E., Zheng, Z., Yengo, L., Lloyd-Jones, L.R., Sidorenko, J., Wu, Y., et al.; eQTL-Gen Consortium (2018). Genome-wide association analyses identify 143 risk variants and putative regulatory mechanisms for type 2 diabetes. *Nat. Commun.* 9, 2941.
3. Markunas, C.A., Johnson, E.O., and Hancock, D.B. (2017). Comprehensive evaluation of disease- and trait-specific enrichment for eight functional elements among GWAS-identified variants. *Hum. Genet.* 136, 911–919.
4. Visscher, P.M., Wray, N.R., Zhang, Q., Sklar, P., McCarthy, M.I., Brown, M.A., and Yang, J. (2017). 10 Years of GWAS Discovery: Biology, Function, and Translation. *Am. J. Hum. Genet.* 101, 5–22.
5. Roman, T.S., and Mohlke, K.L. (2018). Functional genomics and assays of regulatory activity detect mechanisms at loci for lipid traits and coronary artery disease. *Curr. Opin. Genet. Dev.* 50, 52–59.
6. Lu, Q., Powles, R.L., Abdallah, S., Ou, D., Wang, Q., Hu, Y., Lu, Y., Liu, W., Li, B., Mukherjee, S., et al. (2017). Systematic tissue-specific functional annotation of the human genome highlights immune-related DNA elements for late-onset Alzheimer's disease. *PLoS Genet.* 13, e1006933.
7. Claussnitzer, M., Cho, J.H., Collins, R., Cox, N.J., Dermitzakis, E.T., Hurles, M.E., Kathiresan, S., Kenny, E.E., Lindgren, C.M., MacArthur, D.G., et al. (2020). A brief history of human disease genetics. *Nature* 577, 179–189.
8. Jugessur, A., Farlie, P.G., and Kilpatrick, N. (2009). The genetics of isolated orofacial clefts: From genotypes to subphenotypes. *Oral Dis* 15, 437–453.
9. Mangold, E., Ludwig, K.U., and Nöthen, M.M. (2011). Breakthroughs in the genetics of orofacial clefting. *Trends Mol. Med.* 17, 725–733.
10. Mossey, P.A., and Modell, B. (2012). Epidemiology of oral clefts 2012: An international perspective. *Front. Oral Biol* 16, 1–18.
11. Grosen, D., Bille, C., Pedersen, J.K., Skytthe, A., Murray, J.C., and Christensen, K. (2010). Recurrence risk for offspring of twins discordant for oral cleft: A population-based cohort study of the Danish 1936-2004 cleft twin cohort. *Am. J. Med. Genet. A* 152A, 2468–2474.
12. Christensen, K., Juel, K., Herskind, A.M., and Murray, J.C. (2004). Long term follow up study of survival associated with cleft lip and palate at birth. *BMJ* 328, 1405.
13. Yu, Y., Zuo, X., He, M., Gao, J., Fu, Y., Qin, C., Meng, L., Wang, W., Song, Y., Cheng, Y., et al. (2017). Genome-wide analyses of non-syndromic cleft lip with palate identify 14 novel loci and genetic heterogeneity. *Nat. Commun.* 8, 14364.
14. van Rooij, I.A., Ludwig, K.U., Welzenbach, J., Ishorst, N., Thonissen, M., Galesloot, T.E., Ongkosuwito, E., Bergé, S.J., Aldhore, K., Rojas-Martinez, A., et al. (2019). Non-syndromic cleft lip with or without cleft palate: Genome-wide association study in Europeans identifies a suggestive risk locus at 16p12.1 and supports SH3PXD2A as a clefting susceptibility gene. *Genes (Basel)* 10, 1023.
15. Moreno, L.M., Mansilla, M.A., Bullard, S.A., Cooper, M.E., Busch, T.D., Machida, J., Johnson, M.K., Brauer, D., Krahn, K., Daack-Hirsch, S., et al. (2009). FOXE1 association with both isolated cleft lip with or without cleft palate, and isolated cleft palate. *Hum. Mol. Genet.* 18, 4879–4896.

16. Mostowska, A., Gaczowska, A., Żukowski, K., Ludwig, K.U., Hozyasz, K.K., Wójcicki, P., Mangold, E., Böhmer, A.C., Heilmann-Heimbach, S., Knapp, M., et al. (2018). Common variants in DLG1 locus are associated with non-syndromic cleft lip with or without cleft palate. *Clin. Genet.* *93*, 784–793.
17. Leslie, E.J., Carlson, J.C., Shaffer, J.R., Feingold, E., Wehby, G., Laurie, C.A., Jain, D., Laurie, C.C., Doheny, K.F., McHenry, T., et al. (2016). A multi-ethnic genome-wide association study identifies novel loci for non-syndromic cleft lip with or without cleft palate on 2p24.2, 17q23 and 19q13. *Hum. Mol. Genet.* *25*, 2862–2872.
18. Mangold, E., Ludwig, K.U., Birnbaum, S., Baluardo, C., Ferrian, M., Herms, S., Reutter, H., de Assis, N.A., Chawa, T.A., Mattheisen, M., et al. (2010). Genome-wide association study identifies two susceptibility loci for nonsyndromic cleft lip with or without cleft palate. *Nat. Genet.* *42*, 24–26.
19. Ludwig, K.U., Mangold, E., Herms, S., Nowak, S., Reutter, H., Paul, A., Becker, J., Herberz, R., AlChawa, T., Nasser, E., et al. (2012). Genome-wide meta-analyses of nonsyndromic cleft lip with or without cleft palate identify six new risk loci. *Nat. Genet.* *44*, 968–971.
20. Beaty, T.H., Murray, J.C., Marazita, M.L., Munger, R.G., Ruczinski, I., Hetmanski, J.B., Liang, K.Y., Wu, T., Murray, T., Fallin, M.D., et al. (2010). A genome-wide association study of cleft lip with and without cleft palate identifies risk variants near MAFB and ABCA4. *Nat. Genet.* *42*, 525–529.
21. Ludwig, K.U., Böhmer, A.C., Bowes, J., Nikolić, M., Ishorst, N., Wyatt, N., Hammond, N.L., Gözl, L., Thieme, F., Barth, S., et al. (2017). Imputation of orofacial clefting data identifies novel risk loci and sheds light on the genetic background of cleft lip ± cleft palate and cleft palate only. *Hum. Mol. Genet.* *26*, 829–842.
22. Sun, Y., Huang, Y., Yin, A., Pan, Y., Wang, Y., Wang, C., Du, Y., Wang, M., Lan, F., Hu, Z., et al. (2015). Genome-wide association study identifies a new susceptibility locus for cleft lip with or without a cleft palate. *Nat. Commun.* *6*, 6414.
23. Ludwig, K.U., Ahmed, S.T., Böhmer, A.C., Sangani, N.B., Varghese, S., Klamt, J., Schuenke, H., Gültepe, P., Hofmann, A., Rubini, M., et al. (2016). Meta-analysis Reveals Genome-Wide Significance at 15q13 for Nonsyndromic Clefting of Both the Lip and the Palate, and Functional Analyses Implicate GREM1 As a Plausible Causative Gene. *PLoS Genet.* *12*, e1005914.
24. Beaty, T.H., Taub, M.A., Scott, A.F., Murray, J.C., Marazita, M.L., Schwender, H., Parker, M.M., Hetmanski, J.B., Balakrishnan, P., Mansilla, M.A., et al. (2013). Confirming genes influencing risk to cleft lip with/without cleft palate in a case-parent trio study. *Hum. Genet.* *132*, 771–781.
25. Zuccherro, T.M., Cooper, M.E., Maher, B.S., Daack-Hirsch, S., Nepomuceno, B., Ribeiro, L., Caprau, D., Christensen, K., Suzuki, Y., Machida, J., et al. (2004). Interferon Regulatory Factor 6 (IRF6) Gene Variants and the Risk of Isolated Cleft Lip or Palate. *N. Engl. J. Med.* *351*, 769–780.
26. Leslie, E.J., Carlson, J.C., Shaffer, J.R., Butali, A., Buxó, C.J., Castilla, E.E., Christensen, K., Deleyiannis, F.W.B., Leigh Field, L., Hecht, J.T., et al. (2017). Genome-wide meta-analyses of nonsyndromic orofacial clefts identify novel associations between FOXE1 and all orofacial clefts, and TP63 and cleft lip with or without cleft palate. *Hum. Genet.* *136*, 275–286.
27. Birnbaum, S., Ludwig, K.U., Reutter, H., Herms, S., Steffens, M., Rubini, M., Baluardo, C., Ferrian, M., Almeida de Assis, N., Alblas, M.A., et al. (2009). Key susceptibility locus for non-syndromic cleft lip with or without cleft palate on chromosome 8q24. *Nat. Genet.* *41*, 473–477.
28. Rahimov, F., Marazita, M.L., Visel, A., Cooper, M.E., Hitchler, M.J., Rubini, M., Domann, F.E., Govil, M., Christensen, K., Bille, C., et al. (2008). Disruption of an AP-2 α binding site in an IRF6 enhancer is associated with cleft lip. *Nat. Genet.* *40*, 1341–1347.
29. Thieme, F., and Ludwig, K.U. (2017). The Role of Noncoding Genetic Variation in Isolated Orofacial Clefts. *J. Dent. Res.* *96*, 1238–1247.
30. Attanasio, C., Nord, A.S., Zhu, Y., Blow, M.J., Li, Z., Liberton, D.K., Morrison, H., Plajzer-Frick, I., Holt, A., Hosseini, R., et al. (2013). Fine tuning of craniofacial morphology by distant-acting enhancers. *Science* *342*, 1241006.
31. Uslu, V.V., Petretich, M., Ruf, S., Langenfeld, K., Fonseca, N.A., Marioni, J.C., and Spitz, F. (2014). Long-range enhancers regulating Myc expression are required for normal facial morphogenesis. *Nat. Genet.* *46*, 753–758.
32. Liu, H., Leslie, E.J., Carlson, J.C., Beaty, T.H., Marazita, M.L., Lidral, A.C., and Cornell, R.A. (2017). Identification of common non-coding variants at 1p22 that are functional for non-syndromic orofacial clefting. *Nat. Commun.* *8*, 14759.
33. Dunham, I., Kundaje, A., Aldred, S.F., Collins, P.J., Davis, C.A., Doyle, F., Epstein, C.B., Frietze, S., Harrow, J., Kaul, R., et al.; ENCODE Project Consortium (2012). An integrated encyclopedia of DNA elements in the human genome. *Nature* *489*, 57–74.
34. Kundaje, A., Meuleman, W., Ernst, J., Bilenky, M., Yen, A., Heravi-Moussavi, A., Kheradpour, P., Zhang, Z., Wang, J., Ziller, M.J., et al.; Roadmap Epigenomics Consortium (2015). Integrative analysis of 111 reference human epigenomes. *Nature* *518*, 317–330.
35. Aguet, F., Brown, A.A., Castel, S.E., Davis, J.R., He, Y., Jo, B., Mohammadi, P., Park, Y.S., Parsana, P., Segrè, A.V., et al. (2017). Genetic effects on gene expression across human tissues. *Nature* *550*, 204–213.
36. Leslie, E.J., Mansilla, M.A., Biggs, L.C., Schuette, K., Bullard, S., Cooper, M., Dunnwald, M., Lidral, A.C., Marazita, M.L., Beaty, T.H., et al. (2012). Expression and mutation analyses implicate ARHGAP29 as the etiologic gene for the cleft lip with or without cleft palate locus identified by genome-wide association on chromosome 1p22. *Birth Defects Res. A Clin. Mol. Teratol.* *94*, 934–942.
37. Rada-Iglesias, A., Bajpai, R., Prescott, S., Brugmann, S.A., Swigut, T., and Wysocka, J. (2012). Epigenomic annotation of enhancers predicts transcriptional regulators of human neural crest. *Cell Stem Cell* *11*, 633–648.
38. Prescott, S.L., Srinivasan, R., Marchetto, M.C., Grishina, I., Narvaiza, I., Selleri, L., Gage, F.H., Swigut, T., and Wysocka, J. (2015). Enhancer divergence and cis-regulatory evolution in the human and chimp neural crest. *Cell* *163*, 68–83.
39. Wilderman, A., VanOudenhove, J., Kron, J., Noonan, J.P., and Cotney, J. (2018). High-Resolution Epigenomic Atlas of Human Embryonic Craniofacial Development. *Cell Rep* *23*, 1581–1597.
40. Dixon, J.R., Selvaraj, S., Yue, F., Kim, A., Li, Y., Shen, Y., Hu, M., Liu, J.S., and Ren, B. (2012). Topological domains in mammalian genomes identified by analysis of chromatin interactions. *Nature* *485*, 376–380.
41. Freire-Pritchett, P., Schoenfelder, S., Várnai, C., Wingett, S.W., Cairns, J., Collier, A.J., García-Vílchez, R., Furlan-Magaril, M., Osborne, C.S., Fraser, P., et al. (2017). Global reorganisation

- of *cis*-regulatory units upon lineage commitment of human embryonic stem cells. *eLife* 6, e21926.
42. Marchini, J., Howie, B., Myers, S., McVean, G., and Donnelly, P. (2007). A new multipoint method for genome-wide association studies by imputation of genotypes. *Nat. Genet.* 39, 906–913.
 43. Spielman, R.S., McGinnis, R.E., and Ewens, W.J. (1993). Transmission test for linkage disequilibrium: the insulin gene region and insulin-dependent diabetes mellitus (IDDM). *Am. J. Hum. Genet.* 52, 506–516.
 44. Willer, C.J., Li, Y., and Abecasis, G.R. (2010). METAL: fast and efficient meta-analysis of genomewide association scans. *Bioinformatics* 26, 2190–2191.
 45. Bulik-Sullivan, B.K., Loh, P.R., Finucane, H.K., Ripke, S., Yang, J., Patterson, N., Daly, M.J., Price, A.L., Neale, B.M., Corvin, A., et al.; Schizophrenia Working Group of the Psychiatric Genomics Consortium (2015). LD Score regression distinguishes confounding from polygenicity in genome-wide association studies. *Nat. Genet.* 47, 291–295.
 46. de Leeuw, C.A., Mooij, J.M., Heskes, T., and Posthuma, D. (2015). MAGMA: Generalized Gene-Set Analysis of GWAS Data. *PLoS Comput. Biol.* 11, e1004219.
 47. Watanabe, K., Stringer, S., Frei, O., Umičević Mirkov, M., de Leeuw, C., Polderman, T.J.C., van der Sluis, S., Andreassen, O.A., Neale, B.M., and Posthuma, D. (2019). A global overview of pleiotropy and genetic architecture in complex traits. *Nat. Genet.* 51, 1339–1348.
 48. Ernst, J., and Kellis, M. (2015). Large-scale imputation of epigenomic datasets for systematic annotation of diverse human tissues. *Nat. Biotechnol.* 33, 364–376.
 49. Ernst, J., and Kellis, M. (2017). Chromatin-state discovery and genome annotation with ChromHMM. *Nat. Protoc.* 12, 2478–2492.
 50. Lausch, M., Bartusel, M., Alirzayeva, H., Karaolidou, A., Rehim, R., Crispatzu, G., Nikolic, M., Bleckwehl, T., Kolovos, P., van Ijcken, W.F.J., et al. (2018). Disruption of the TFAP2A Regulatory Domain Causes Banchio-Oculo-Facial Syndrome (BOFS) and Illuminates Pathomechanisms for Other Human Neurocristopathies. *Cell Stem Cell* 24, 736–752, e12.
 51. Schmidt, E.M., Zhang, J., Zhou, W., Chen, J., Mohlke, K.L., Chen, Y.E., and Willer, C.J. (2015). GREGOR: evaluating global enrichment of trait-associated variants in epigenomic features using a systematic, data-driven approach. *Bioinformatics* 31, 2601–2606.
 52. Mukherjee, K., Ishii, K., Pillalammarri, V., Kammin, T., Atkin, J.F., Hickey, S.E., Xi, Q.J., Zepeda, C.J., Gusella, J.F., Talkowski, M.E., et al. (2016). Actin capping protein CAPZB regulates cell morphology, differentiation, and neural crest migration in craniofacial morphogenesis. *Hum. Mol. Genet.* 25, 1255–1270.
 53. McLennan, R., Bailey, C.M., Schumacher, L.J., Teddy, J.M., Morrison, J.A., Kasemeier-Kulesa, J.C., Wolfe, L.A., Gogol, M.M., Baker, R.E., Maini, P.K., et al. (2017). DAN (NBL1) promotes collective neural crest migration by restraining uncontrolled invasion. *J. Cell Biol.* 16, 3339–3354.
 54. Matt, N., Dupé, V., Garnier, J.M., Dennefeld, C., Chambon, P., Mark, M., and Ghyselinck, N.B. (2005). Retinoic acid-dependent eye morphogenesis is orchestrated by neural crest cells. *Development* 132, 4789–4800.
 55. Kousa, Y.A., Fuller, E., and Schutte, B.C. (2018). IRF6 and AP2A Interaction Regulates Epidermal Development. *J. Invest. Dermatol.* 138, 2578–2588.
 56. Lin-Shiao, E., Lan, Y., Welzenbach, J., Alexander, K.A., Zhang, Z., Knapp, M., Mangold, E., Sammons, M., Ludwig, K.U., and Berger, S.L. (2019). p63 establishes epithelial enhancers at critical craniofacial development genes. *Sci. Adv.* 5, eaaw0946.
 57. Kobayashi, K., Jakt, L.M., and Nishikawa, S.I. (2013). Epigenetic regulation of the neuroblastoma genes, *Arid3b* and *Mycn*. *Oncogene* 32, 2640–2648.
 58. Takebe, A., Era, T., Okada, M., Martin Jakt, L., Kuroda, Y., and Nishikawa, S. (2006). Microarray analysis of PDGFR α + populations in ES cell differentiation culture identifies genes involved in differentiation of mesoderm and mesenchyme including ARID3b that is essential for development of embryonic mesenchymal cells. *Dev. Biol.* 293, 25–37.
 59. Jakobsen, L.P., Borup, R., Vestergaard, J., Larsen, L.A., Lage, K., Maroun, L.L., Kjaer, I., Niemann, C.U., Andersen, M., Knudsen, M.A., et al. (2009). Expression analyses of human cleft palate tissue suggest a role for osteopontin and immune related factors in palatal development. *Exp. Mol. Med.* 41, 77–85.
 60. Leslie, E.J., Carlson, J.C., Shaffer, J.R., Buxó, C.J., Castilla, E.E., Christensen, K., Deleyiannis, F.W.B., Field, L.L., Hecht, J.T., Moreno, L., et al. (2017). Association studies of low-frequency coding variants in nonsyndromic cleft lip with or without cleft palate. *Am. J. Med. Genet. A* 173, 1531–1538.
 61. Cox, L.L., Cox, T.C., Moreno Uribe, L.M., Zhu, Y., Richter, C.T., Nidey, N., Standley, J.M., Deng, M., Blue, E., Chong, J.X., et al. (2018). Mutations in the Epithelial Cadherin-p120-Catenin Complex Cause Mendelian Non-Syndromic Cleft Lip with or without Cleft Palate. *Am. J. Hum. Genet.* 102, 1143–1157.
 62. Marini, N.J., Asrani, K., Yang, W., Rine, J., and Shaw, G.M. (2019). Accumulation of rare coding variants in genes implicated in risk of human cleft lip with or without cleft palate. *Am. J. Med. Genet. A* 179, 1260–1269.
 63. Lupiáñez, D.G., Spielmann, M., and Mundlos, S. (2016). Breaking TADs: How Alterations of Chromatin Domains Result in Disease. *Trends Genet.* 32, 225–237.
 64. Lupiáñez, D.G., Kraft, K., Heinrich, V., Krawitz, P., Brancati, F., Klopocki, E., Horn, D., Kayserili, H., Opitz, J.M., Laxova, R., et al. (2015). Disruptions of topological chromatin domains cause pathogenic rewiring of gene-enhancer interactions. *Cell* 161, 1012–1025.
 65. Lausch, M., Bartusel, M., Rehim, R., Alirzayeva, H., Karaolidou, A., Crispatzu, G., Zentis, P., Nikolic, M., Bleckwehl, T., Kolovos, P., et al. (2019). Modeling the Pathological Long-Range Regulatory Effects of Human Structural Variation with Patient-Specific hiPSCs. *Cell Stem Cell* 24, 736–752.e12.
 66. Vermillion, K.L., Lidberg, K.A., and Gammill, L.S. (2014). Expression of actin-binding proteins and requirement for actin-depolymerizing factor in chick neural crest cells. *Dev. Dyn.* 243, 730–738.
 67. Ernst, J., and Kellis, M. (2012). ChromHMM: automating chromatin-state discovery and characterization. *Nat. Methods* 9, 215–216.
 68. Richardson, R.J., Hammond, N.L., Coulombe, P.A., Saloranta, C., Nousiainen, H.O., Salonen, R., Berry, A., Hanley, N., Headon, D., Karikoski, R., and Dixon, M.J. (2014). Periderm prevents pathological epithelial adhesions during embryogenesis. *J. Clin. Invest.* 124, 3891–3900.
 69. Fraser, H.B., Lam, L.L., Neumann, S.M., and Kobor, M.S. (2012). Population-specificity of human DNA methylation. *Genome Biol.* 13, R8.

70. Liu, J., Hutchison, K., Perrone-Bizzozero, N., Morgan, M., Sui, J., and Calhoun, V. (2010). Identification of genetic and epigenetic marks involved in population structure. *PLoS ONE* 5, e13209.
71. Butali, A., Mossey, P.A., Adeyemo, W.L., Eshete, M.A., Gowans, L.J.J., Busch, T.D., Jain, D., Yu, W., Huan, L., Laurie, C.A., et al. (2019). Genomic analyses in African populations identify novel risk loci for cleft palate. *Hum. Mol. Genet.* 28, 1038–1051.
72. Mukhopadhyay, N., Bishop, M., Mortillo, M., Chopra, P., Hetmanski, J.B., Taub, M.A., Moreno, L.M., Valencia-Ramirez, L.C., Restrepo, C., Wehby, G.L., et al. (2020). Whole genome sequencing of orofacial cleft trios from the Gabriella Miller Kids First Pediatric Research Consortium identifies a new locus on chromosome 21. *Hum. Genet* 139, 215–226.
73. Carlson, J.C., Anand, D., Butali, A., Buxo, C.J., Christensen, K., Deleyiannis, F., Hecht, J.T., Moreno, L.M., Orioli, I.M., Padilla, C., et al. (2019). A systematic genetic analysis and visualization of phenotypic heterogeneity among orofacial cleft GWAS signals. *Genet. Epidemiol.* 43, 704–716.
74. Huang, L., Jia, Z., Shi, Y., Du, Q., Shi, J., Wang, Z., Mou, Y., Wang, Q., Zhang, B., Wang, Q., et al. (2019). Genetic factors define CPO and CLO subtypes of nonsyndromic orofacial cleft. *PLoS Genet.* 15, e1008357.
75. Ray, A.K., Marazita, M.L., Pathak, R., Beever, C.L., Cooper, M.E., Goldstein, T., Shaw, D.F., and Field, L.L. (2004). TP63 mutation and clefting modifier genes in an EEC syndrome family. *Clin. Genet.* 66, 217–222.
76. Liu, H., Busch, T., Eliason, S., Anand, D., Bullard, S., Gowans, L.J.J., Nidey, N., Petrin, A., Augustine-Akpan, E.A., Saadi, I., et al. (2017). Exome sequencing provides additional evidence for the involvement of ARHGAP29 in Mendelian orofacial clefting and extends the phenotypic spectrum to isolated cleft palate. *Birth Defects Res.* 109, 27–37.
77. Savastano, C.P., Brito, L.A., Faria, Á.C., Setó-Salvia, N., Peskett, E., Musso, C.M., Alvizi, L., Ezquina, S.A.M., James, C., GOS-gene, et al. (2017). Impact of rare variants in ARHGAP29 to the etiology of oral clefts: role of loss-of-function vs missense variants. *Clin. Genet.* 91, 683–689.
78. Mansouri, A. (1998). The role of Pax3 and Pax7 in development and cancer. *Crit. Rev. Oncog.* 9, 141–149.
79. Bishop, M.R., Diaz Perez, K.K., Sun, M., Ho, S., Chopra, P., Mukhopadhyay, N., Hetmanski, J.B., Taub, M.A., Moreno-Urbe, L.M., Valencia-Ramirez, L.C., et al. (2020). Genome-wide Enrichment of De Novo Coding Mutations in Orofacial Cleft Trios. *Am. J. Hum. Genet.* 107, 124–136.

HGGA, Volume 2

Supplemental information

**Integrative approaches generate insights
into the architecture of non-syndromic cleft lip
with or without cleft palate**

Julia Welzenbach, Nigel L. Hammond, Miloš Nikolić, Frederic Thieme, Nina Ishorst, Elizabeth J. Leslie, Seth M. Weinberg, Terri H. Beaty, Mary L. Marazita, Elisabeth Mangold, Michael Knapp, Justin Cotney, Alvaro Rada-Iglesias, Michael J. Dixon, and Kerstin U. Ludwig

Supplemental Material and Methods

Cohort Description

The meta-analysis included data from three previously published individual GWAS on nsCL/P (Bonn case-control GWAS cohort¹, GENEVA trio cohort², POFC GWAS cohort³, Table S1). We included all nsCL/P summary statistics that were publicly accessible until June 2018. In the present study we combined the three GWAS cohorts to generate the largest nsCL/P meta-analysis to date. In accordance with previous studies⁴⁻⁶, two meta-analyses were performed: (1) using all individuals with diverse population backgrounds (MAiC, for **Meta-Analysis in Clefting**) and (2) using the European data sets only (MAiC_{Euro}). Data quality control (QC) included the detection and removal of overlapping individuals, confirmation of ethnicity, and data re-analysis as described in the following:

Data QC on the individual Bonn and GENEVA cohorts was done as previously described^{4,5}. For the POFC GWAS cohort, imputed genotypes were retrieved from dbGaP. Briefly, this data comprised genotypes for 557,677 single nucleotide polymorphisms (SNPs) in 11,855 individuals of diverse ethnicities (Ethiopia, Nigeria, China, India, Philippines, Denmark, Hungary, Spain, Turkey, Argentina, Colombia, Guatemala, Puerto Rico, United States), representing 3,981 families. Based on these genotype data, a dataset was constructed that had maximal overlap to the original published POFC GWAS³ while excluding individuals that have already been analysed previously as part of the GENEVA cohort. First we generated a combined POFC-GENEVA dataset and used KING⁷ on a set of 115,380 genotyped variants, to estimate relatedness between individuals. Relationship was defined based on a KING kinship coefficient ≥ 0.0884 (representing 2nd degree relationship), and affected individuals or families were removed from the POFC cohort. In the remaining POFC individuals, we then aimed at maximizing the number of case-parent trios, resulting in 1,328 complete trios (1,319 had been included in Leslie et al. 2016³; nine additional nsCL/P trios were identified based on inference of family structure). This data set formed the final POFC_{trio} cohort. From the remaining families (where no case-parent trio had been drawn from), independent individuals were selected to construct the POFC_{case-control} cohort.

Individuals were designated “cases” if affected with nsCL/P, based on the phenotypic data provided. In situations when multiple individuals within one family were affected, the individual with the lowest number of missing genotypes was included. This resulted in 848 cases. For the control set, data from 1,568 families without any affected individual were available. From these families, the individual with lowest number of missing genotypes was selected as control and included in the POFC_{case-control} cohort (n=1,568 controls).

Ethnicity was identified based on the information provided in column “country of origin”, and genotype data. Individuals were classified into “European” and “Non-European” based on the smallest distance in mean and standard deviation of the first two principal component analysis (PCA)-Eigenvectors, to the defined Central European country category ‘Denmark/Hungary/Spain/Turkey’. To exclude individuals that were identical within studies (or were part of a “superfamily” with other individuals in the cohorts), the relationship between individuals was calculated with KING based on the shared variants in the Bonn, GENEVA and POFC cohorts. Again, individuals showing a kinship-coefficient of ≥ 0.0884 were excluded, this resulted in removal of 90 case-parent trios, one case, and seven controls from the POFC cohort. Thus, the final MAiC dataset resulted in 848 cases, 1,561 controls and 1,238 case-parent trios after sample QC (Table S1).

Statistical analyses

Statistical analyses were performed separately for case-control cohorts and case-parent trios, respectively. Notably, the Ludwig et al. (2017) meta-analysis had applied the FBAT-dosage method to generate genotypes for the trio cohorts, lacking individual relative risk (RR) information. In the present study, best-guess genotypes were assigned based on *a-posteriori* genotype probabilities of ≥ 0.6 . In the case-control cohorts, GWAS was performed using SNPTEST and -method expected, by incorporating 5-18 dimensions of the multi-dimensional-scaling coordinates⁸, respectively. For the case-parent trios a transmission disequilibrium test (TDT) was performed on the best-guess genotypes⁹. Given the present sample sizes, we accounted for the limited power of imputation approaches to correctly predict rare and low-frequency variants by retaining robust SNPs only (info-score ≥ 0.4 , minor allele frequency

(MAF) >1 % in controls and non-transmitting parents). Moreover, for the case-parent trios, SNPs had to be present in ≥75 % of the families. After data cleaning procedures, we meta-analyzed the GWAS data of all four sub-cohorts (Bonn case-control, GENEVA case-parent trios, POFC case-control and POFC case-parent trios) using METAL¹⁰. METAL combines p-values across studies while considering directions of effects and effective sub-cohort size, as indicated by N_{eff} . Here, N_{eff} is defined as

$$case/control N_{eff} = \frac{4}{\frac{1}{n_{cases}} + \frac{1}{n_{controls}}} \quad \text{and} \quad trio N_{eff} = \frac{4}{\frac{1}{n_{trio}} + \frac{1}{n_{trio}}}$$

Post-analysis SNPs that were absent from more than one sub-cohort (Bonn case-control, GENEVA case-parent trios, POFC case-control and POFC case-parent trios) were removed. Thus, the final MAiC dataset contained 7,744,527 SNPs in MAiC, and 7,690,843 SNPs in MAiC_{Euro}. To estimate the SNP-based heritability (h^2) for nsCL/P on the liability scale, we generated a European case-control-only dataset that excluded the case-parent trios (Table S1). Using this data set we performed linkage disequilibrium (LD) score regression as implemented in *ldsc*¹¹, for individuals of European ethnicity. Because LD score regression requires a homogeneous data structure, it was not possible to apply *ldsc* on the whole nsCL/P dataset with mixed ethnicities and complex cohort structure (case/control and case-parent trios).

Identification of novel nsCL/P risk loci

We defined ‘previously identified risk loci for nsCL/P’ as those having reached genome-wide significance ($P < 5 \times 10^{-08}$) in either GWAS, meta-analysis or large-scale systematic study before ($n=40$, Table S2). For each of the lead variants (as defined in the respective original study) a window of strong LD ($r^2 \geq 0.6$) was defined. The most distant SNPs at this threshold determined the boundaries for known genetic risk loci. For those of the previously identified risk loci that reached genome-wide significance in the MAiC dataset ($n=26$), the 1000 Genome Phase3 reference panel, Central European backbone, was used to compute r^2 in the FUMA (Functional Mapping and Annotation of Genome-Wide Association Studies) platform, v1.3.2¹². For the remaining 14 loci, $r^2 = 0.6$ boundaries were estimated using LDproxy provided in LDlink 3.2.0 suite¹³, in a mixed (European/East Asian) population. ‘Novel risk loci’ were defined as those with $P < 5 \times 10^{-08}$ in the MAiC analyses if located outside of the previously

identified loci. For each of the novel risk loci identified in MAiC, we defined the associated region using the same parameters (Table S2). For follow-up analyses, recent data from GTEx (v8)¹⁴ were assessed using the GTEx browser.

Gene-based and pathway analyses

Gene-based analyses in MAiC and MAiC_{Euro} were performed using MAGMA¹⁵ (v1.06), implemented in FUMA. MAGMA's gene analysis uses a multiple regression approach to properly incorporate LD between markers and to detect multi-marker effects. We run the gene-based analysis with default parameters (SNP-wide mean model), using 1000 Genome Phase3 as reference panel and the full distribution of imputed input SNPs of MAiC (N=7,744,527; info-score \geq 0.4; MAF $>$ 1%). These were mapped to 17,911 protein-coding genes based to a distance of 0kb upstream/downstream of the genes, resulting in threshold of test-wide significance of $P = 2.79 \times 10^{-6}$ (i.e., 0.05/17,911).

Epigenetic datasets for mid-facial development

Identification of datasets relevant to mid-facial development

Human cell-type and developmental-stage specific data for mid-facial development is underrepresented (or not represented at all) in large consortia data such as ENCODE¹⁶. However, available data in the Gene Expression Omnibus (GEO) covered mid-facial development from (i) early stages (hNCC¹⁷, accessed through GSE28874), (ii) differentiated human cNCC¹⁸ (accessed through GSE70751), and (iii) embryonic craniofacial human tissue of different CS (accessed through GSE97752)¹⁹. Each of the datasets is briefly described in the following section:

(i) *Human neural crest cells*. For the establishment of hNCCs, an *in vitro* differentiation model had been used in which hESC (H9 cell line) were first induced to form neuroectodermal spheres (hNEC) that subsequently gave rise to migratory cells expressing early NC markers and recapitulating neuronal, mesenchymal and melanocytic differentiation potential of the neural crest^{17,20}. Sequential ChIP assays

in hNCCs had been performed from approximately 10^7 hNCC cells per experiment²¹, as described before²² and included chromatin modifications H3K27ac, H3K4me1, H3K4me3 and H3K27me3.

(ii) *Human cranial neural crest cells*. Human cNCCs had been differentiated *in vitro* from iPSCs (H9 cell line), first forming hNECs which then later attached and gave rise to migratory cNCCs which could be maintained up to 18 passages¹⁸. In cNCC, ChIPs had been performed using approximately 0.5×10^{-7} to 1.0×10^{-7} cells per experiment, with the same protocol as described for hNCC. ChIP-Seq had been performed for chromatin modifications H3K27ac, H3K4me1, H3K4me3 and H3K27me3, as for hNCC.

(iii) *Craniofacial tissue from different stages*. Human embryonic CT was collected, staged, and provided by the Joint MRC/Wellcome Trust Human Developmental Biology Resource (HDBR, www.hdbr.org). Information describing the developmental stage, termination method, collection site, and karyotype of each embryo and the ChIP protocol is provided in the original study¹⁹. Briefly, samples encompassed CS 13 (4.5 weeks *post-conception*, wpc, 5 embryos), 14 (5 wpc, 3 embryos), 15 (5.5 wpc, 3 embryos), 17 (6 wpc, 4 embryos) and 20 (8 wpc, one embryo), and one sample of 10 wpc. For samples of CS13-17, ChIP-Seq was performed for chromatin modifications H3K27ac, H3K4me1, H3K4me3, H3K27me3 and H3K36me3 and resulted in a mean of 37.3 million uniquely aligned reads per sample and chromatin mark (Table S4). Data for H3K9me3 marks were imputed. For CS20 and 10 wpc, H3K27ac3 ChIP-Seq data was experimentally derived, all other marks were imputed (Table S3).

ChIP-Seq Data processing

For hNCC and cNCC, raw data were available in fastq format. ChIP-seq data from craniofacial data Wilderman et al. (2018)¹⁹ comprise processed formats, including imputed signals, peaks and segmentation data. In order to ensure comparability between the three data sources, computational processing of ChIP-seq data as published in Wilderman et al. 2018 (QC, alignment, peak calling, epigenetic imputation, chromatin segmentation) was adopted to the hNCC/cNCC bioinformatics pipeline. Scripts have been used as deposited on <https://github.com/cotneylab/ChIP-Seq>.

Briefly, FastQC (v0.11.7) was used to combine multiple fastq files of one ChIP experiment and perform QC. Alignment of the single-end reads with length of 36 bp (hNCC) and 36-50 bp (cNCC) to the human genome (hg19) was performed with Bowtie2 (v2.3.2)²³. Fragment sizes of each library were estimated using PhantomPeakQualTools (v1.14)²⁴. For hNCC and cNCC, the number of uniquely aligned reads per sample and chromatin mark is given in Table S5. In the following, analysis of treatment against input sample was performed to generate p-value based signal tracks and peak files based on estimated library fragment size using MACS2 (v2.1.1.20160309)²⁵.

Chromatin imputation and segmentation.

To obtain uniform data sets, chromatin imputation followed by chromatin state segmentation was performed. First, H3K9me3 and H3K36me3 marks in both hNCC and cNCC were imputed using ChromImpute (v1.0.1)²⁶, based on 127 cell types from the Roadmap Epigenome Project²⁷. Briefly we used *P*-value-based signal files for marks H3K27ac, H3K4me1, H3K4me3, H3K27me3, H3K9me3 and H3K36me3 for 127 tissues and cell types. Conversion from bigWig to bedGraph format was done using the ENCODE function 'BigWigToBedGraph'. Both the hNCC and cNCC p-value signals in bedGraph format as well as the Roadmap reference samples were converted to 25 bp resolution and processed for model training and generation of imputed signals.

Imputed hNCC and cNCC signal files for each individual chromosome and each chromatin mark were binarized, and segmentation was performed using the core+K27ac 18-state chromatin model provided by Roadmap with ChromHMM²⁸ which uses a multivariate Hidden Markov Model that explicitly models the combinatorial presence or absence of each mark to predict 18 chromatin states. This procedure can identify tissue-specific regulatory information from initial tissue/cells for which no gene expression data is available. Because of the low number of chromatin marks measured in the NCC samples, epigenetic imputation issues and the higher risk of batch effect between hNCC, cNCC and CT, we adopted a robust strategy and condensed the 18 generated states into eight states.

Translation of genetic associations into tissue- and time point-specific regulatory effects at a systematic level

Enrichment analyses using GREGOR

Based on chromatin segments obtained from hNCC, cNCC and CT (see section “Chromatin imputation and segmentation”), we used the GREGOR software (Genomic Regulatory Elements and Gwas Overlap algorithm)²⁹ to evaluate the enrichment of significant SNPs from the MAiC data in the available regulatory features (i.e., eight predicted chromatin states). A set of samples from the Roadmap Epigenomics project (comprising both fetal and adult tissue samples) was selected as independent dataset for comparison. These included ESC H1 cell line (ESC.H1), iPS cell line (IPSC.DF.19.11), bone marrow derived cultured mesenchymal stem cells (STRM.MRW.MSC), primary B cells from peripheral blood (BLD.CD19.PPC), foreskin fibroblast primary cells skin01 (SKIN.PEN.FRISK.FIB.01), fetal muscle trunk (MUS.TRNK.FET), fetal muscle leg (MUS.LEG.FET), fetal stomach (GI.STMC.FET), psoas muscle (MUS.PSOAS), rectal mucosa donor 29 (GI.RECT.MUC.29), and HeLa-S3 cervical carcinoma cell line (CRVX.HELAS3.CNCR). As input we used MAiC nsCL/P variants with $P \leq 0.001$ without additional variants in LD ($N=22,999$); this threshold was selected to balance between adequate statistical power and true positive association signals. To test if the observed enrichment is specific for nsCL/P, we configured a control SNP set comprising an equal number of SNPs with $P > 0.1$ from MAiC data, which were selected to represent the same MAF distribution as the test input.

Characterization of nsCL/P risk variants and candidate gene prioritization in context of epigenetic mid-facial time line.

For comprehensive insights in regulatory mechanisms at nsCL/P risk loci, we finally integrated all available genetic and functional data (MAiC associations, $GWAS_{TAD}$ - and $r^2 \geq 0.6$ -region boundaries, NCC- and CT-specific active chromatin sites, chromatin segmentation tracks and pChi-C *cis* interactions). Based on this approach we attempt to prioritize genetic variants with regulatory effect and potential

downstream target genes and to detect relevant regulatory elements specific for the early (hNCC/cNCC) or later mid-facial development (CT). Chromatin signals that were either only predicted in individual or few cells/tissues were discarded as artefacts. Risk loci without interpretable chromatin signals at the associated region (either no chromatin signals or artefacts) and too complex risk loci (either with many overlapping genes or very broad associated regions that both makes it impossible to derive clear conclusions on regulatory effects on particular genes) were excluded from this part of functional interpretation.

Supplemental Figures

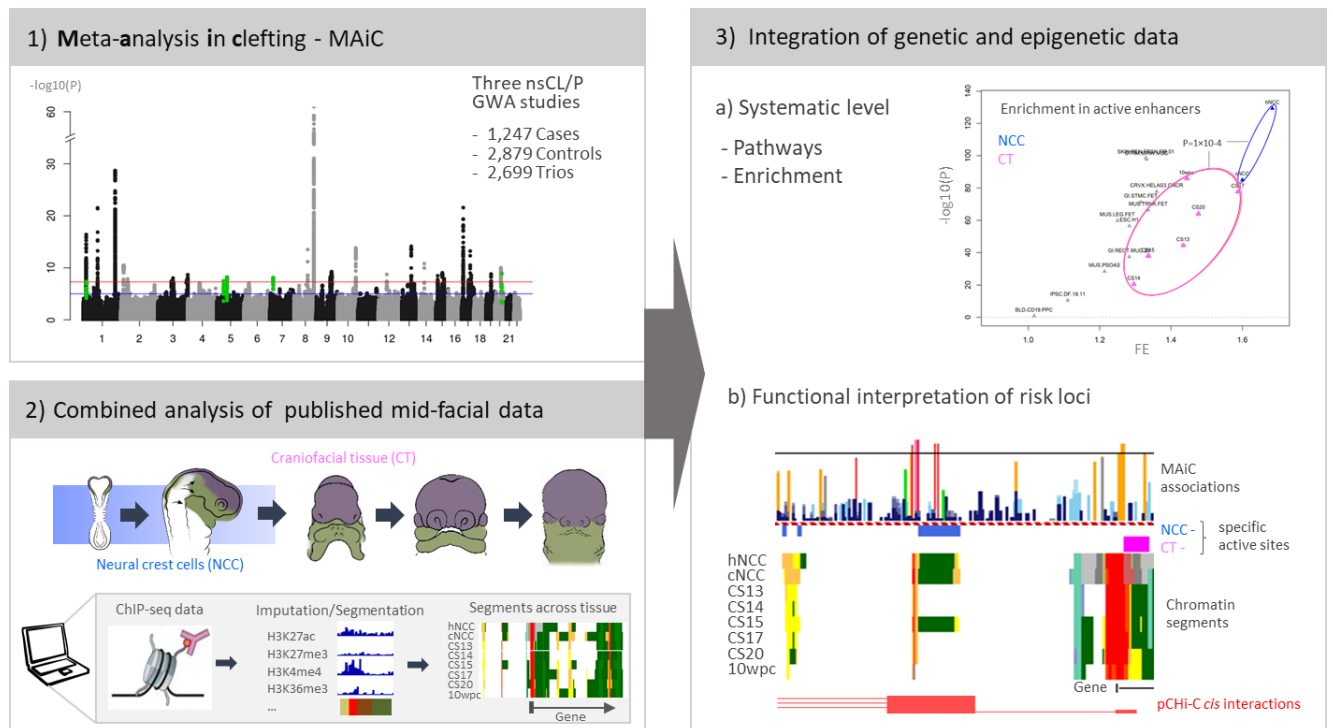


Figure S1 Graphical workflow of the study. **1)** Manhattan plot of summary statistics for MAiC in nonsyndromic cleft lip with/without cleft palate (nsCL/P). Blue and red lines indicate suggestive (at $-\log_{10}(1 \times 10^{-5})$) and genome-wide (at $-\log_{10}(5 \times 10^{-8})$) significance. **2)** Available epigenetic datasets for cells and tissue relevant for craniofacial development were processed in a joint bioinformatic pipeline to generate a comparable map of epigenetic data for the functional analysis of nsCL/P across mid-facial development. **3)** Systematic integration of MAiC association and epigenetic data to analyse the enrichment of genetic variants, to reveal relevant biological pathways and to study regulatory mechanisms at nsCL/P risk loci. GWA - Genome-wide association; ChIP-seq - Chromatin immunoprecipitation followed by sequencing; hNCC - human neural crest cell; cNCC - cranial NCC; CS - Carnegie stage; wpc - weeks *post-conceptum*; pChI-C - Promoter capture chromosome capture.

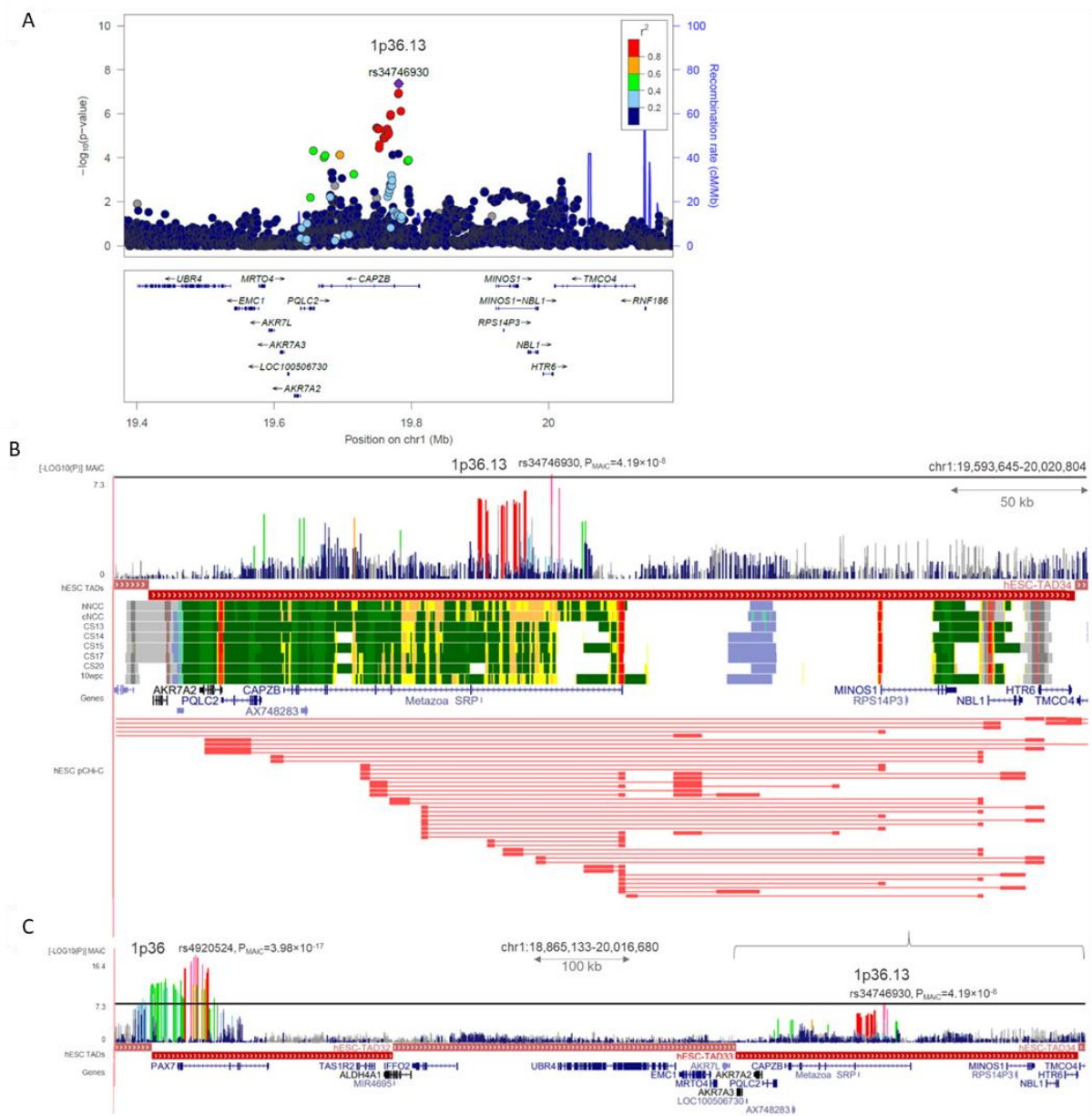


Figure S2 Regulatory architecture at novel risk locus for nonsyndromic cleft lip with/without left palate (nsCL/P) at 1p36.13. A) Regional association plot of chromosomal region 1p36.1 in vicinity of *CAPZB*, with lead variant rs34746930. Data from MAiC, plot generated with LocusZoom. **B)** Zoom into regulatory architecture at 1p36.13. Based on the extent of the topologically associated domain around rs34746930 and variants in linkage disequilibrium (LD) $r^2 \geq 0.6$, different layers of data were aggregated. Tracks include (top-down): MAiC association P-values with colour code based on LD to rs34746930; chromatin segmentation data from early human neural crest cell (hNCC), cranial NCC (cNCC), craniofacial tissue (CT) of different Carnegie stages (CS) and 10 weeks *post-conceptum* (wpc, colour code as in Figure 3 of main text); RefSeq gene positions; and promoter capture Hi-C *cis*-interactions collected in hESC (GSE8682133). **C)** Association structure at 1p36.13, indicating its independence from the known risk locus 1p36 around *PAX7*. Further information on the region is provided as Supplemental Text.

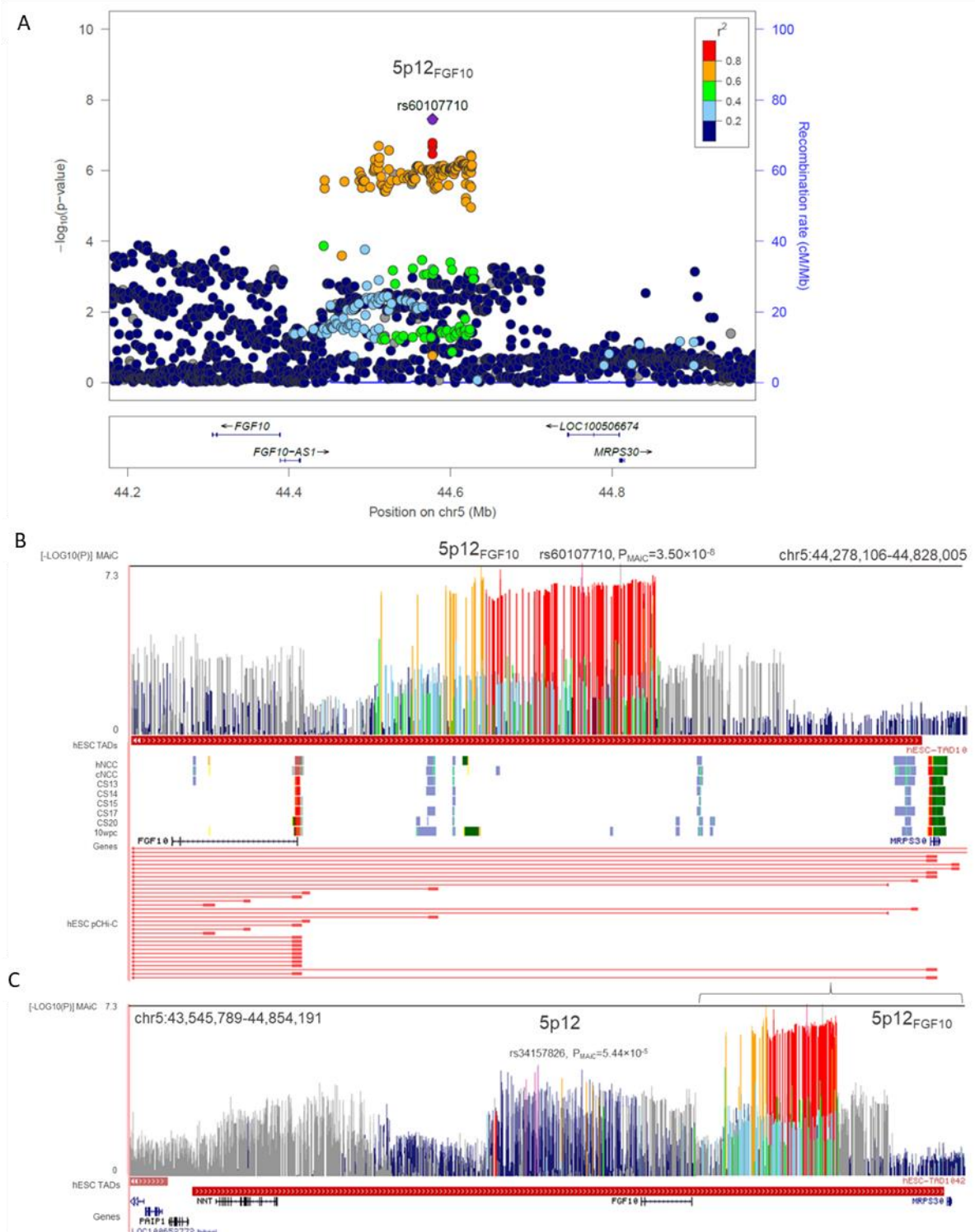


Figure S3 Regulatory architecture at novel risk locus for nonsyndromic cleft lip with/without cleft palate (nsCL/P) risk locus at 5p12_{FGF10}. **A)** Regional association plot of chromosomal region 5p12_{FGF10} with lead variant rs60107710. Data from MAiC, plot generated with LocusZoom. **B)** Zoom into regulatory architecture at 5p12_{FGF10}. Based on the extent of the topologically associated domain around rs60107710 and variants in linkage disequilibrium (LD) $r^2 \geq 0.6$, different layers of data were aggregated. Tracks include (top-down): MAiC association P-values with colour code based on LD to rs60107710; chromatin segmentation data from early human neural crest cell (hNCC), cranial NCC (cNCC), craniofacial tissue (CT) of different Carnegie stages (CS) and 10 weeks *post-conceptum* (wpc, colour code as in Figure 3 of main text); RefSeq gene positions; and promoter capture Hi-C *cis*-interactions collected in hESC (GSE8682133). **C)** Association structure at 5p12_{FGF10}, indicating its independence from the known risk locus 5p12 previously reported in Chinese individuals. Further information on the region is provided as Supplemental Text.

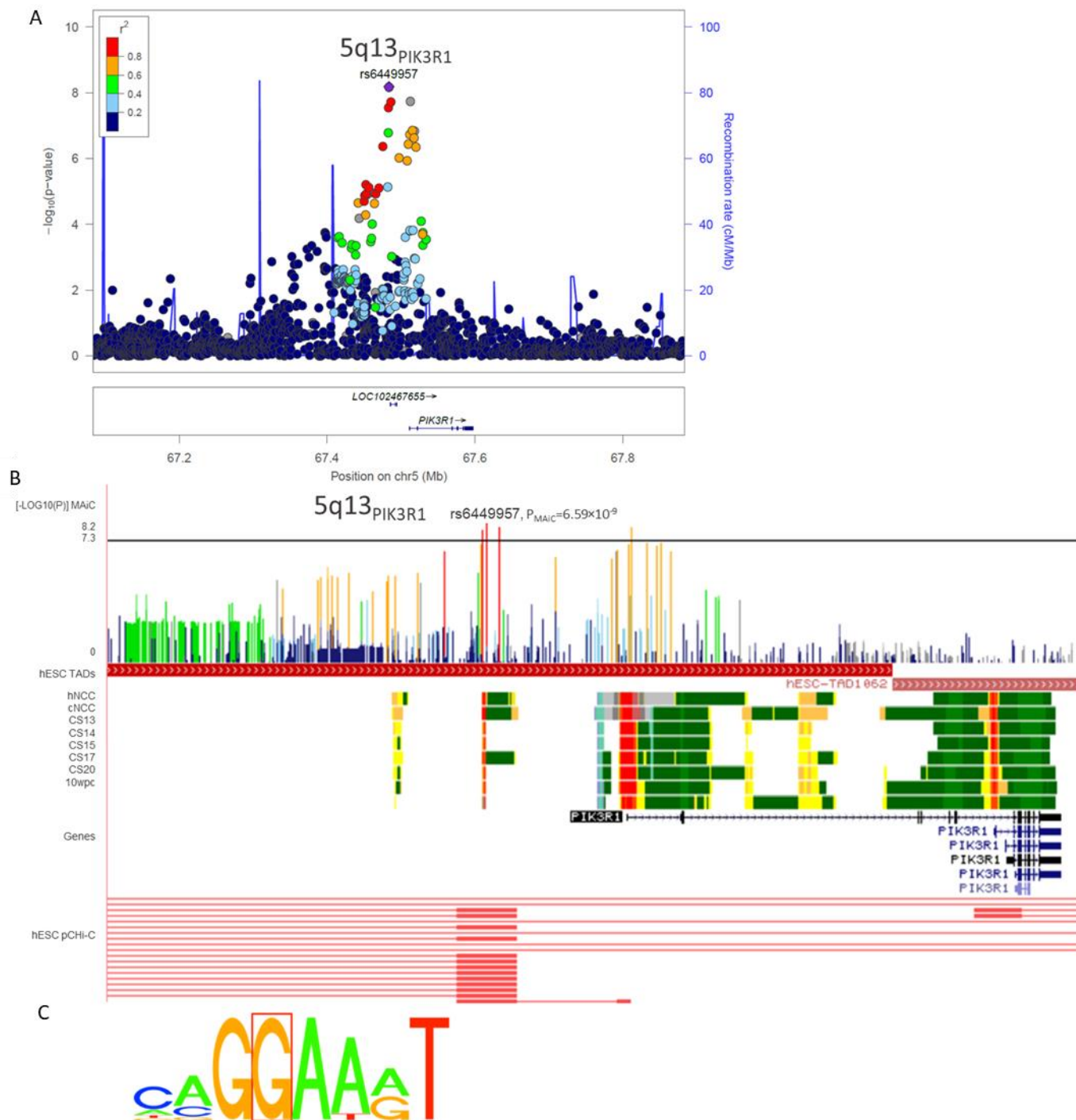


Figure S4 Regulatory architecture at novel risk locus for nonsyndromic cleft lip with/without cleft palate (nsCL/P) risk locus at 5q13_{PIK3R1}. **A)** Regional association plot of chromosomal region 5q13_{PIK3R1} with lead variant rs6449957. Data from MAiC, plot generated with LocusZoom. **B)** Zoom into regulatory architecture at 5q13_{PIK3R1}. Based on the extent of the topologically associated domain around rs6449957 and variants in linkage disequilibrium (LD) $r^2 \geq 0.6$, different layers of data were aggregated. Tracks include (top-down): MAiC association P-values with colour code based on LD to rs6449957; chromatin segmentation data from early human neural crest cell (hNCC), cranial NCC (cNCC), craniofacial tissue (CT) of different Carnegie stages (CS) and 10 weeks *post-conceptum* (wpc, colour code as in Figure 3 of main text); RefSeq gene positions; and promoter capture Hi-C *cis*-interactions collected in hESC (GSE8682133). **C)** According to JASPAR 2018, rs6449956 (G>T), in high LD with rs6449957, is predicted to disrupt a binding motif for FEV, a member of the Ets-family of transcription factors. Further information on the region is provided as Supplemental Text.

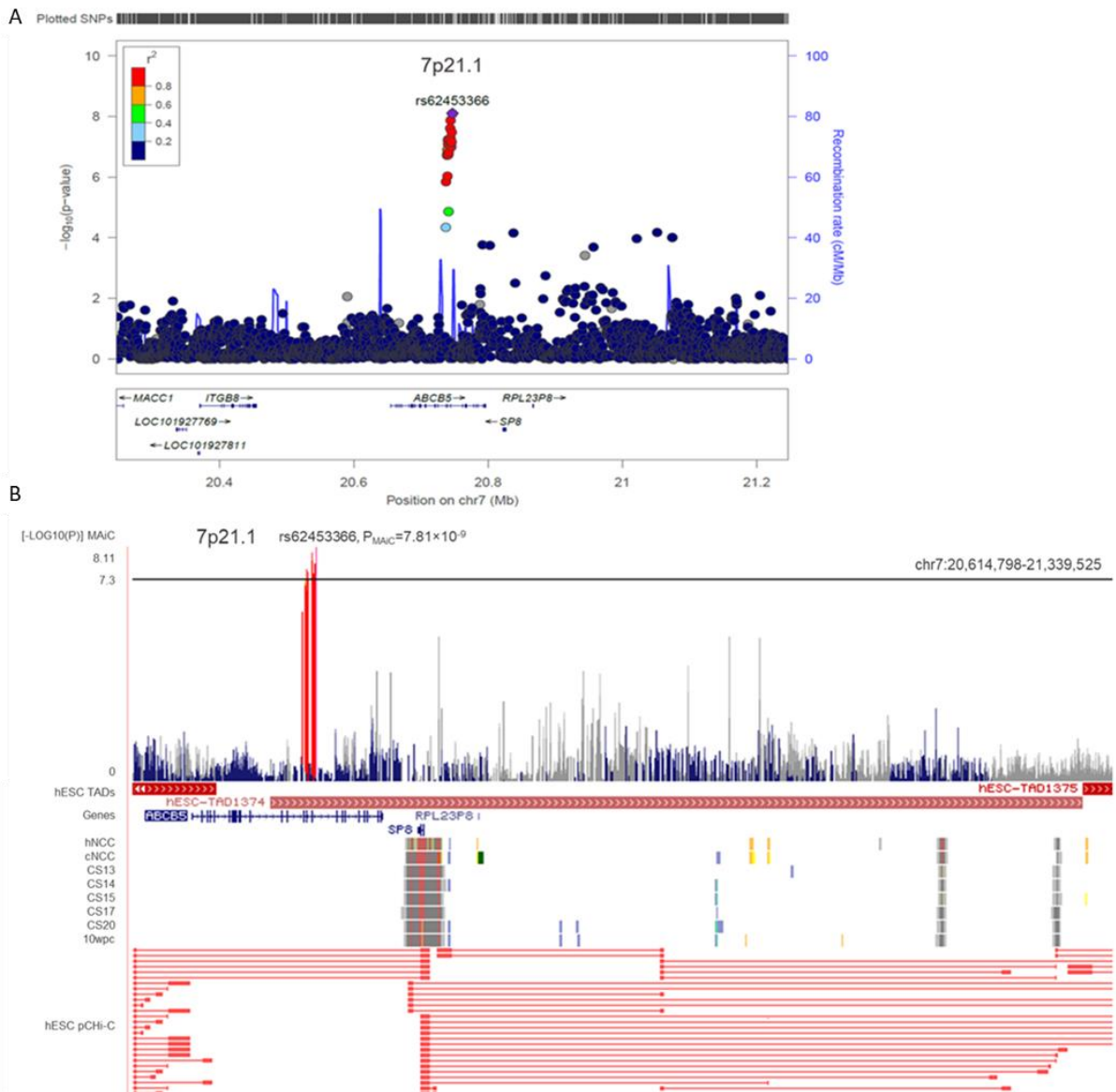


Figure S5 Regulatory architecture at novel risk locus for nonsyndromic cleft lip with/without cleft palate (nsCL/P) risk locus at 7p21.1. **A)** Regional association plot of chromosomal region 7p21.1 with lead variant rs62453366. Data from MAiC, plot generated with LocusZoom. **B)** Zoom into regulatory architecture at 7p21.1. Based on the extent of the topologically associated domain around rs62453366 and variants in linkage disequilibrium (LD) $r^2 \geq 0.6$, different layers of data were aggregated. Tracks include (top-down): MAiC association P-values with colour code based on LD to rs62453366; chromatin segmentation data from early human neural crest cell (hNCC), cranial NCC (cNCC), craniofacial tissue (CT) of different Carnegie stages (CS) and 10 weeks *post-conceptum* (wpc, colour code as in Figure 3 of main text); RefSeq gene positions; and promoter capture Hi-C *cis*-interactions collected in hESC (GSE8682133). Further information on the region is provided as Supplemental Text.

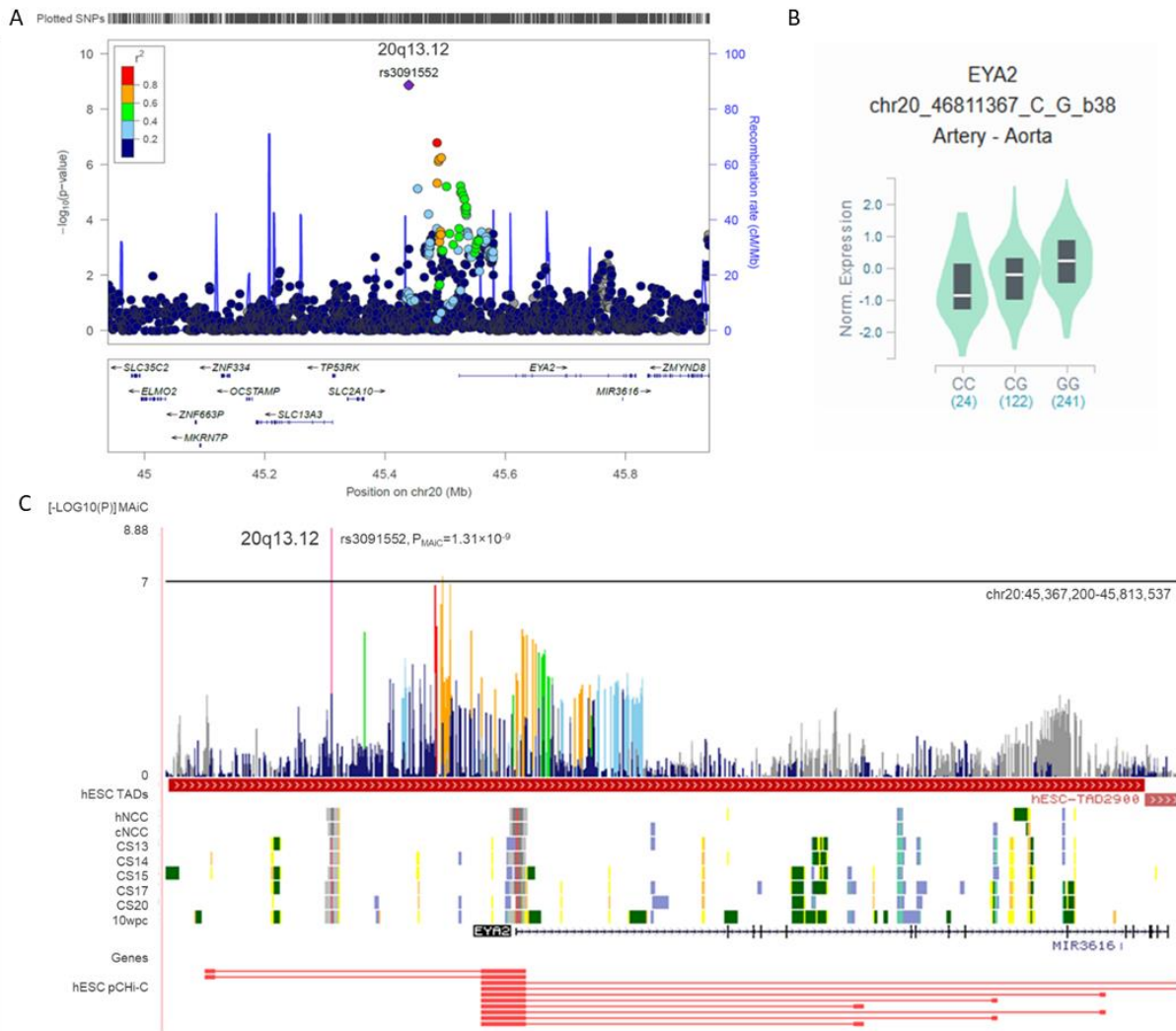


Figure S6 Regulatory architecture at novel risk locus for nonsyndromic cleft lip with/without cleft palate (nsCL/P) risk locus at 20q13.12. A) Regional association plot at chromosomal region 20q13.12, with lead variant rs3091552. Data from MAiC, plot generated with LocusZoom. **B)** According to GTEx(v8) data, rs3091552 is an eQTL for *EYA2* in artery/aorta tissue, the risk allele being G. **C)** Zoom into regulatory architecture at 20q13.12. Based on the extent of the topologically associated domain around rs3091552 and variants in linkage disequilibrium (LD) $r^2 \geq 0.6$, different layers of data were aggregated. Tracks include (top-down): MAiC association P-values with colour code based on LD to rs3091552; chromatin segmentation data from early human neural crest cell (hNCC), cranial NCC (cNCC), craniofacial tissue (CT) of different Carnegie stages (CS) and 10 weeks *post-conceptum* (wpc, colour code as in Figure 3 of main text); RefSeq gene positions; and promoter capture Hi-C *cis*-interactions collected in hESC (GSE8682133). Further information on the region is provided as Supplemental Text.

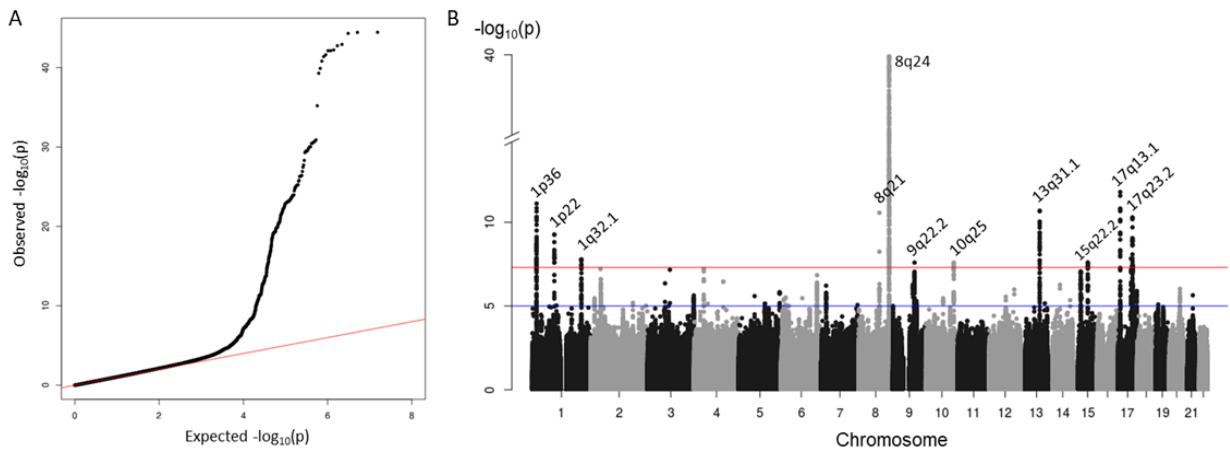


Figure S7 Meta-analysis in clefting (MAiC) performed in the European cohorts. A) In the quantile-quantile-plot, the number and magnitude of observed associations between single variants and nonsyndromic cleft lip with or without cleft palate (nsCL/P) is compared to the association statistics expected under the null hypothesis of no association. The lambda value is 1.04. **B)** In the Manhattan plot, associations of genetic variants across all autosomes and nsCL/P are plotted against the $-\log_{10}$ transformed P-values of MAiC. Level of suggestive significance is highlighted in blue at $-\log_{10}(1 \times 10^{-5})$, and genome-wide significance in red at $-\log_{10}(5 \times 10^{-8})$. N=716 SNPs at ten nsCL/P risk loci were detected as genome-wide significant, the most significant locus being the established 8q24 nsCL/P risk locus (Birbaum et al. 2009). The SNP with the lowest P-value in MAiC_{euro} is rs72728734 (chr8:129,933,720, $P=3.58 \times 10^{-45}$).

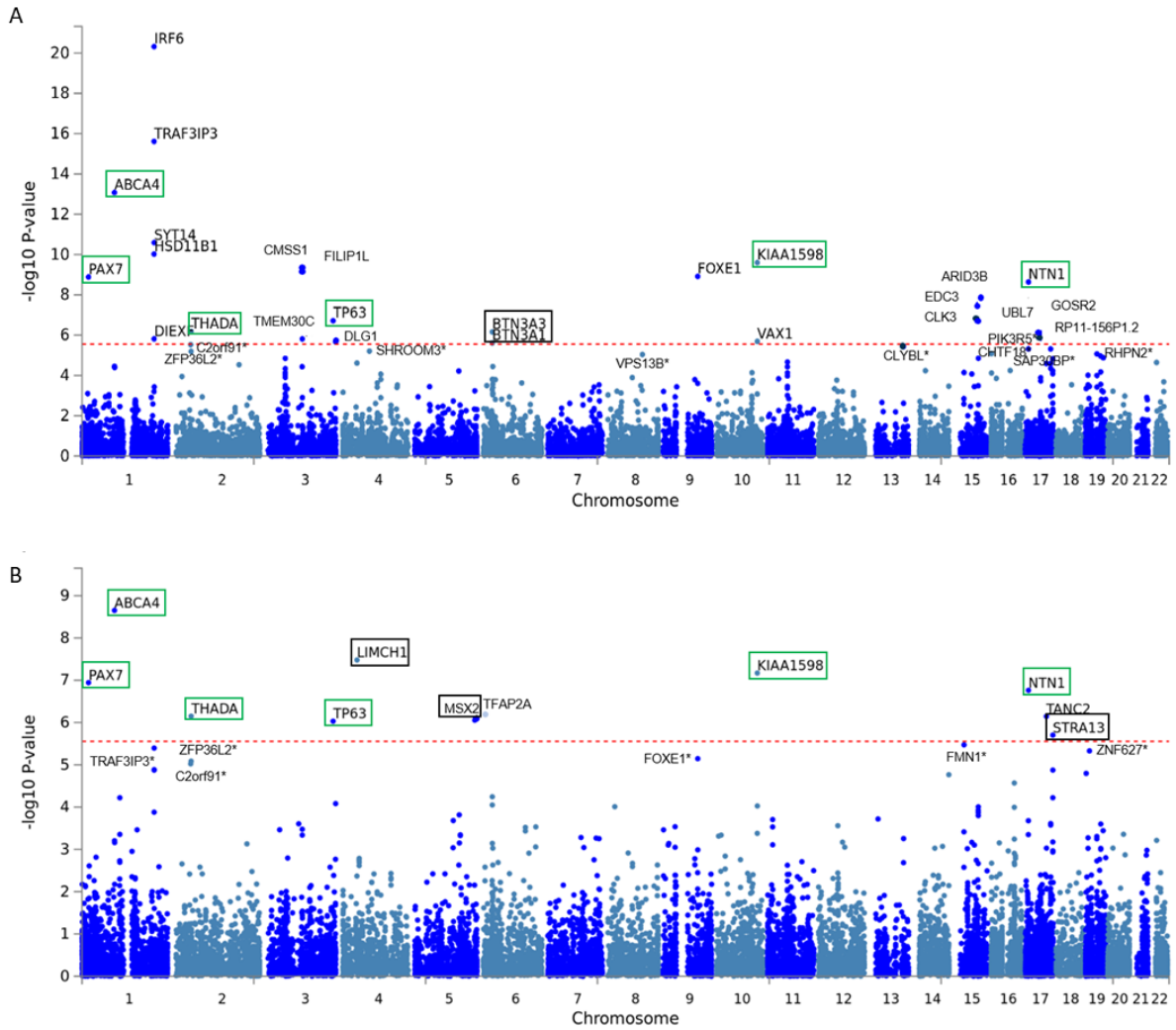


Figure S8 Gene-based Manhattan plot for the multi-ethnic analysis (A) and the European cohort (B). MAiC summary statistics for $n=7,744,527$ (MAiC) / $n=7,690,843$ (MAiC_{Euro}) SNPs were mapped to 17,921 protein coding genes based to a distance of 0kb upstream/downstream of the genes. Of those, nominally significant gene-based associations ($P < 0.05$) were obtained for 1,358 (MAiC) and 1,222 (MAiC_{Euro}) genes, respectively. Test-wide significance (indicated by red dashed line) was defined at $P = 2.79 \times 10^{-6}$ (Online Methods). Twenty-five (MAiC) and eleven (MAiC_{Euro}) genes, respectively, reached test-wide significance. Additionally, nine (MAiC) and six (MAiC_{Euro}) genes were significant at suggestive level ($2.79 \times 10^{-6} < P < 10^{-5}$, highlighted by '*'). Six genes (*ABCA4*, *KIAA1598*, *PAX7*, *NTN1*, *TP63* and *THADA*; highlighted in green boxes) were identified with test-wide significance in both analyses. Across both analyses, five genes were identified at test-wide significance (*BTN3A3*, *BTN3A1*; *LIMCH1*, *MSX2* and *STRA13*; highlighted in black boxes) which do not map to any of the known nsCL/P risk loci. Data generated in FUMA.

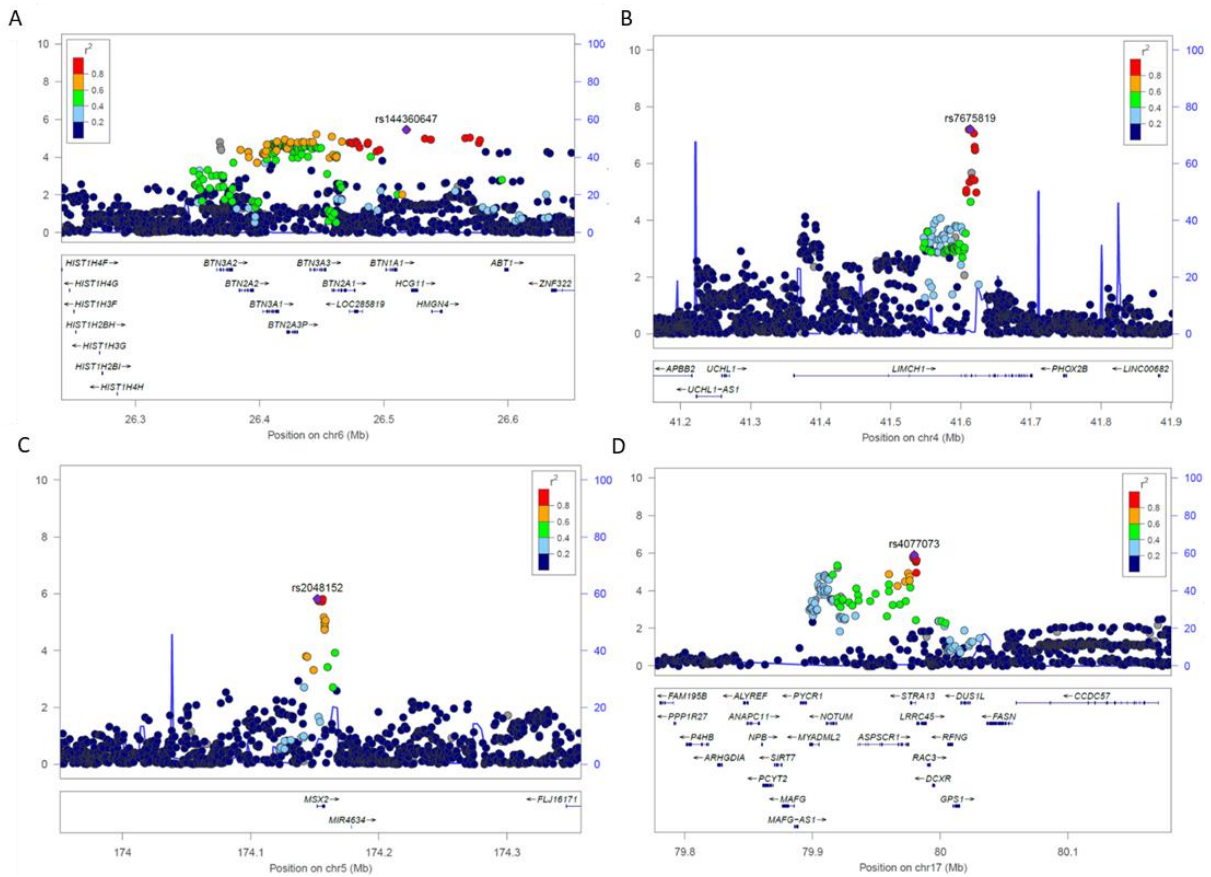


Figure S9 Regional association plots of genes from gene-based analysis which had not been previously reported. Genes were detected in gene-based analysis as implemented in FUMA, applied to MAiC and MAiCEuro summary statistics. Input SNPs were mapped to 17,911 protein coding genes based to a distance of 0kb upstream/downstream. **(A)** Gene-based analysis in MAiC revealed two test-wide significant genes *BTN3A3* ($P=6.96 \times 10^{-7}$) and *BTN3A1* ($P=2.44 \times 10^{-6}$). This region does not map to any of the known nsCL/P risk loci. **(B-D)** In MAiCEuro three test-wide significant genes - *LIMCH1* ($P=3.31 \times 10^{-8}$), *MSX2* ($P=8.80 \times 10^{-7}$) and *STRA13* ($P=1.99 \times 10^{-6}$) were identified outside of the established risk loci for nsCL/P. Plots were generated in LocusZoom.

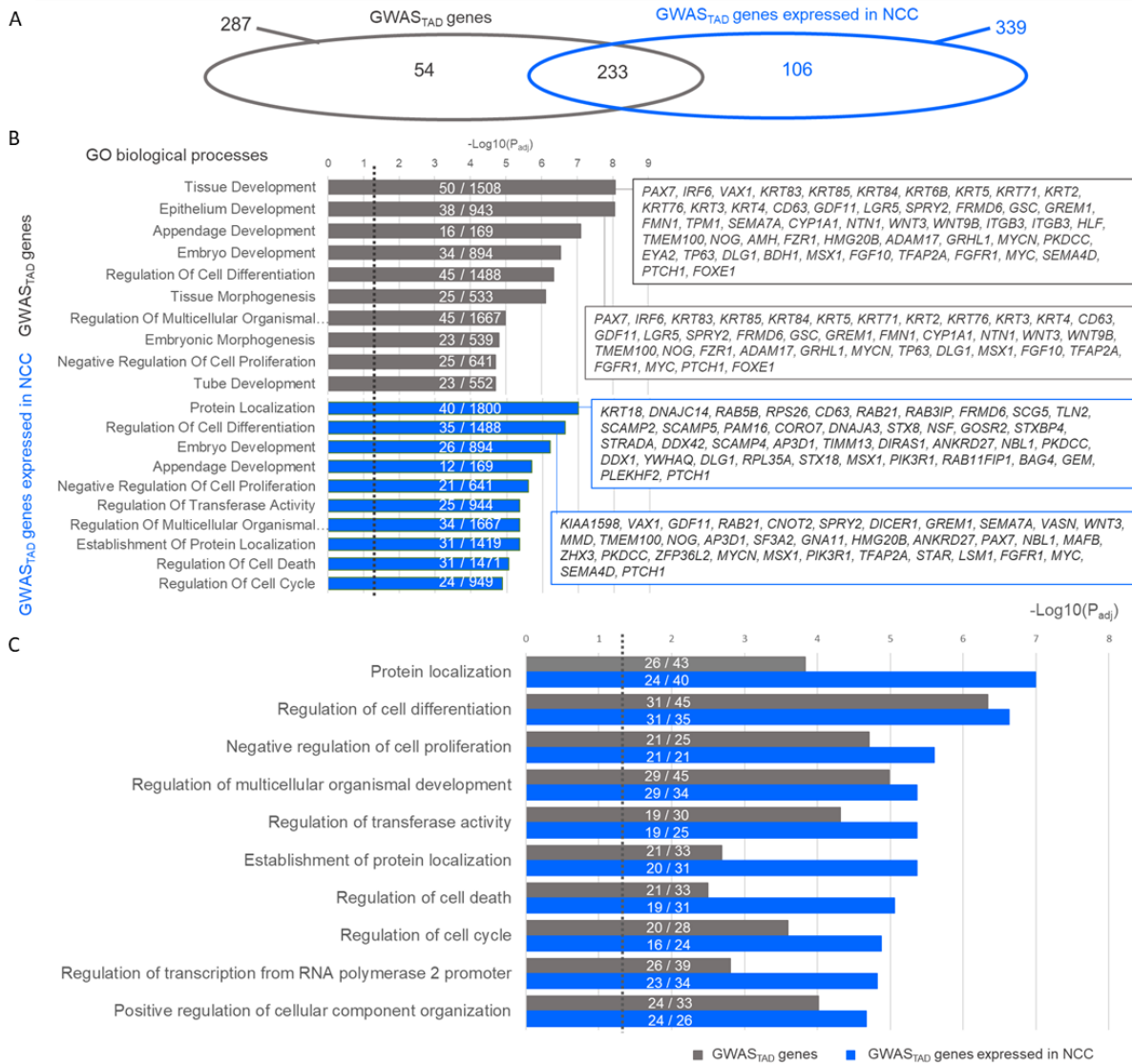


Figure S10 Gene ontology (GO) analysis of 'biological processes' in MAiC data. For each of the 45 nsCL/P risk loci, genes located within the corresponding topological associating domain (TAD) regions were extracted (GWAS_{TAD}-genes). This set of 407 genes was cross-referenced with expression data from neural crest cells (NCC; Laugsch *et al.* (2018) (GSE108522)), revealing expression of 240 GWAS_{TAD} genes (GWAS_{TAD}-genes expressed in NCC). Using the 'GENE2FUNC' application of FUMA (v1.3.4b), enrichment analysis of both gene sets was performed. **A**) Number of significant GO-terms ($P_{adj} \leq 0.05$) in both analyses (n=287 and 339, respectively), with their overlap (n=233) indicated. **B**) Top10 GO biological processes of the individual analyses for 'GWAS_{TAD}' (gray) and 'GWAS_{TAD} genes expressed in NCC' (blue). Dashed line indicates the significance threshold of $P_{adj}=0.05$. Within the bars, the numbers of nsCL/P risk loci / genes represented in this pathway are provided. **C**) Across both analyses, 233 pathways were shared. Of those, 157 had lower P-values in the subset of 'GWAS_{TAD} genes expressed in NCC'. Here, the top10 GO biological processes are shown.

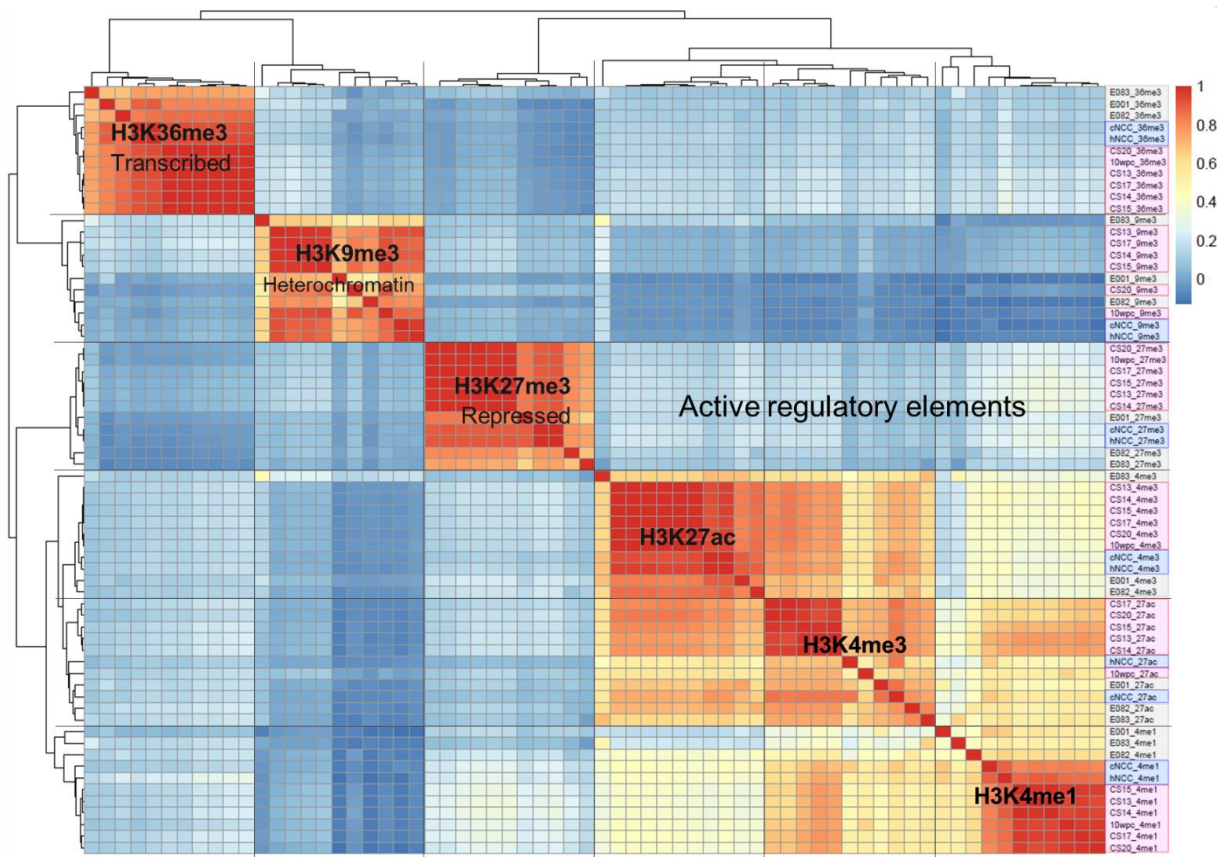


Figure S11 Chromatin modifications in mid-facial development. Hierarchical clustering of pairwise Pearson correlations of epigenetic data. ChIP-seq signals of six histone modifications obtained in human neural crest cells (hNCC); cranial NCC (cNCC, both highlighted in blue); six craniofacial tissue samples (Carnegie stage (CS) 13, CS14, CS15, CS17, CS20, 10 weeks *post-conceptum*, all highlighted in red); and three Roadmap samples (embryonal stem cells (ESC) I3, fetal brain, fetal heart, highlighted in gray).

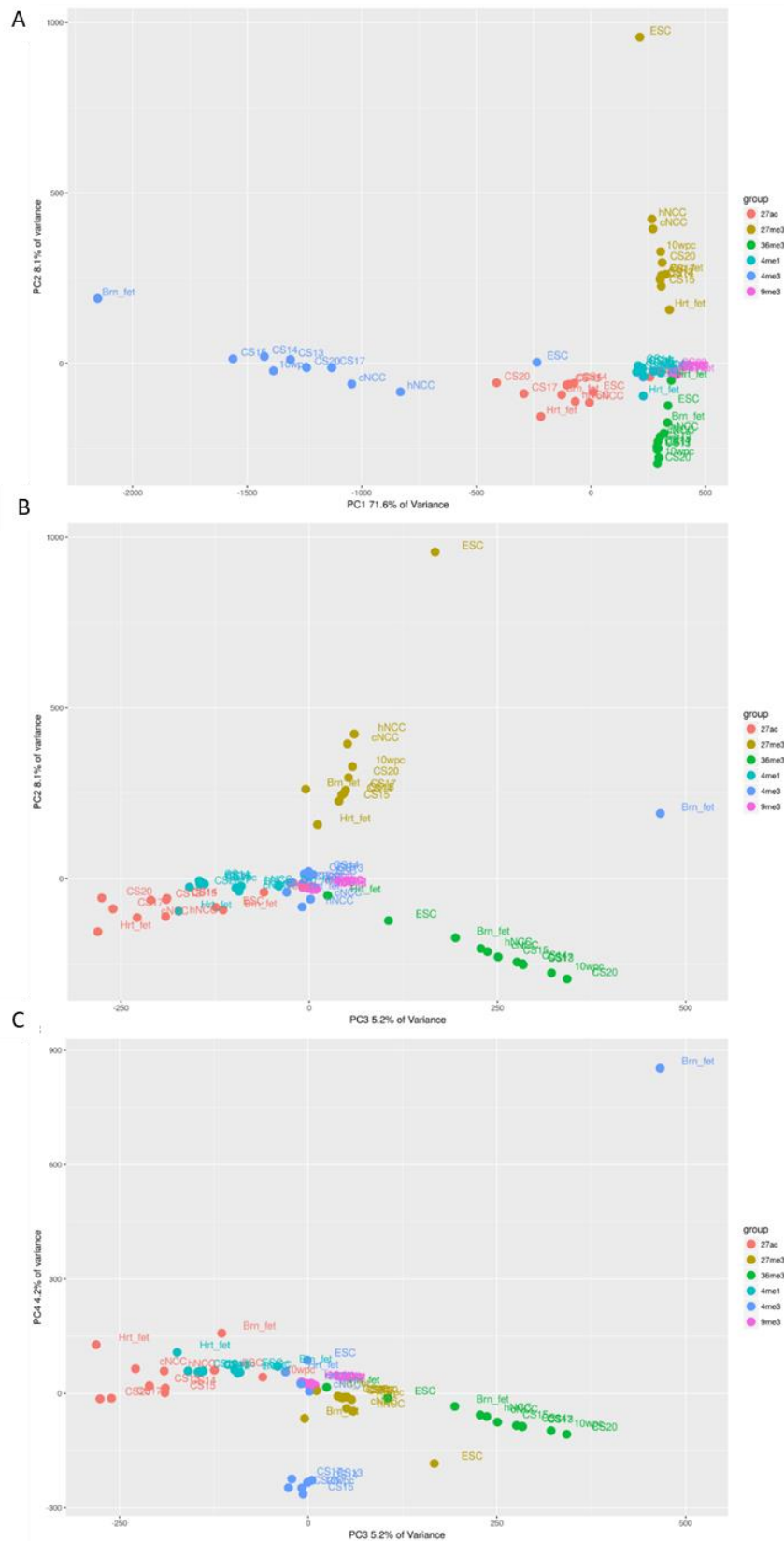


Figure S12 Principal component analysis plot of all imputed chromatin mark signals in neural crest cells (NCC), craniofacial tissue of different Carnegie stages (CS) and selected Roadmap samples. Projection of first vs. second (A), second vs. third (B) and third vs fourth (C) principal component (PC) as analyses based on genome-wide signal profiles of H3K27ac, H3K27me3, H3K36me3, H3K4me1, H3K4me3 and H3K9me3 in early human neural crest cells (hNCC), cranial NCC (cNCC), craniofacial tissue and Roadmap samples of ESC (E001), fetal heart (E083) and fetal brain (E082). Samples are color-coded by chromatin mark. Percentages of variance explained by each PC are indicated along each axis.

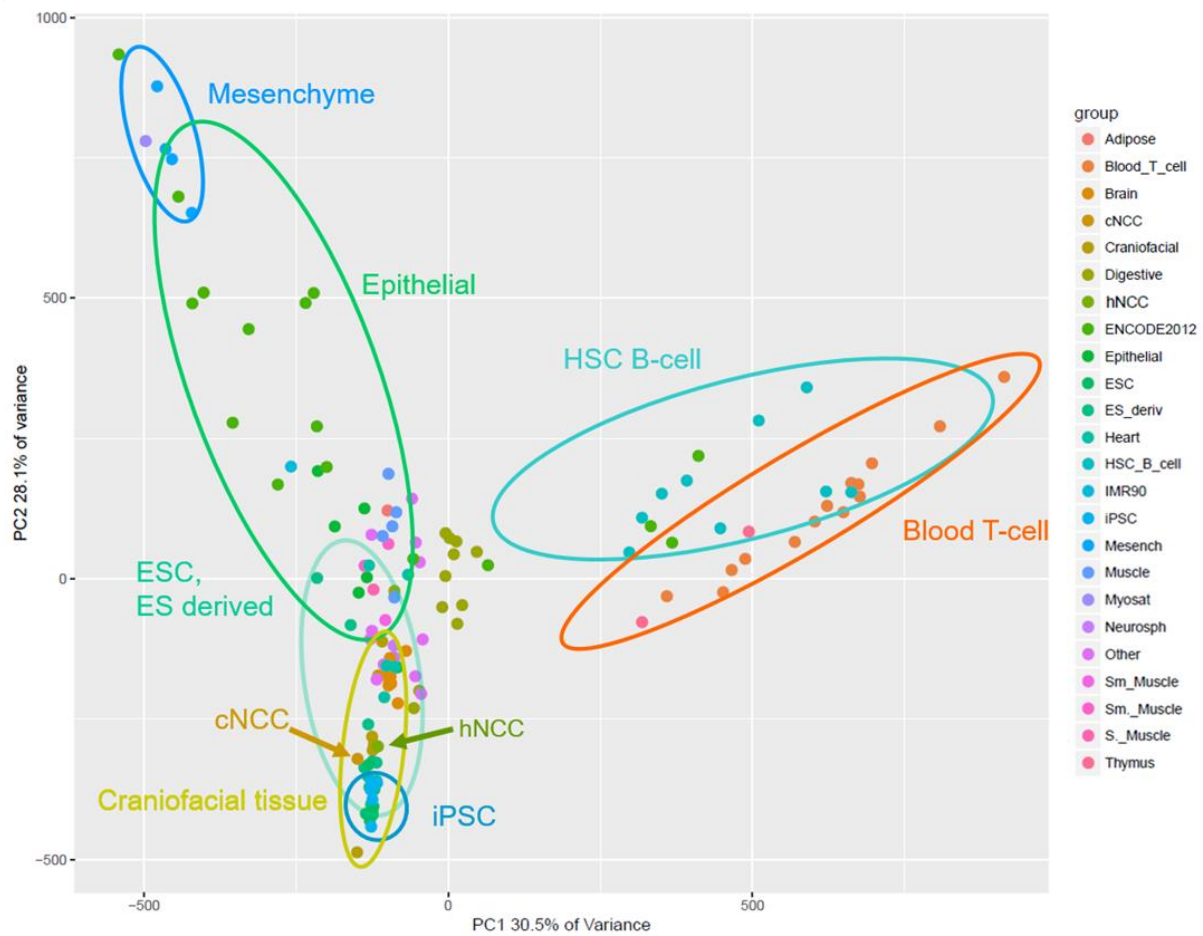


Figure S13 Principal component analysis (PCA) plot based on genome-wide H3K27ac signals in early human neural crest cells (hNCC), cranial NCC (cNCC), craniofacial tissue of different Carnegie stages (CS) all Roadmap samples based on chromatin mark H3K27ac. PCA projection shows the the first and second component dimensions for genome-wide signal profiles of H3K27ac in hNCC, cNCC, craniofacial tissue and all 127 Roadmap/ENCODE samples. Samples are grouped and colour coded as indicated in the legend. Percentages of variance across samples explained by each component are indicated along each axis.

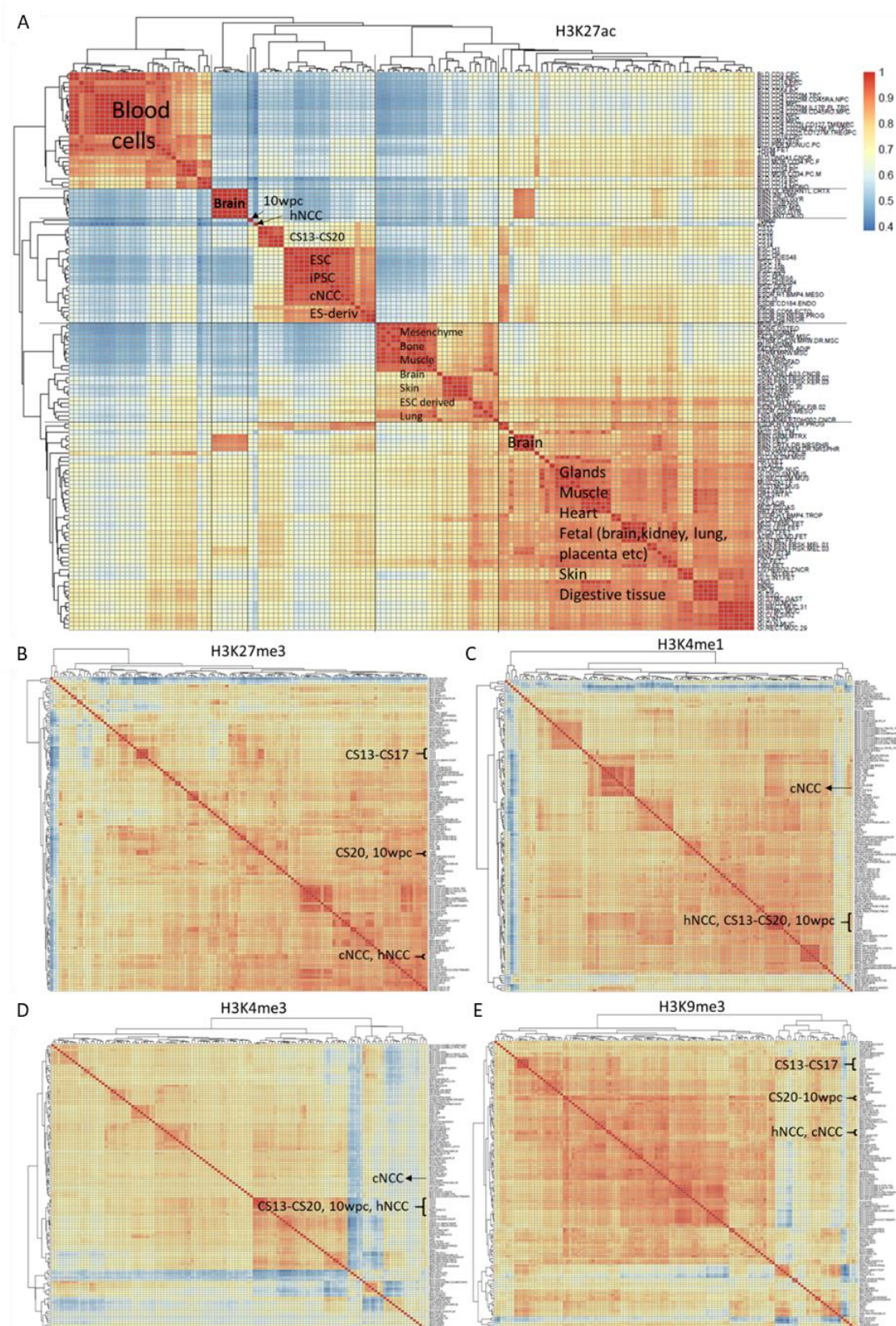


Figure S14 Heatmap and hierarchical clustering of pairwise Pearson correlations. (A) H3K27ac, **(B)** H3K27me3, **(C)** H3K4me1, **(D)** H3K4me3 and **(E)** H3K9me3 signals. For each chromatin modification, the heatmap was generated based on the respective genome-wide ChIP-Seq signals measured in early human neural crest cells (NCC), cranial NCC, six craniofacial tissue samples (Carnegie stage (CS) 13, CS14, CS15, CS17, CS20, 10 weeks *post-conceptum*) and all 127 Roadmap/ENCODE samples. Relatedness of epigenomic profiles by sample is indicated by dendrogram along the axes of the heatmap. Red indicates positive correlation between datasets. The underlying signal comparisons were calculated with deepTools 3.1.3 multiBigwigSummary in bins mode.



Figure S15 Cumulative percentage of chromatin segments on autosomes in neural crest cells (NCC), craniofacial tissue and selected Roadmap samples. Based on the 18-state model (A), data was aggregated into eight states to increase robustness of the analyses (B). For each chromatin state, fractions were calculated in human NCC (hNCC), cranial NCC (cNCC), six craniofacial tissue samples (Carnegie stage (CS) 13, CS14, CS15, CS17, CS20, 10 weeks *post-conceptum*) and a selection of segmentations generated by Roadmap Epigenome (Roadmap Epigenomics Consortium et al., 2015). Color code as presented in the legends, abbreviation of chromatin states as listed in Main Text.

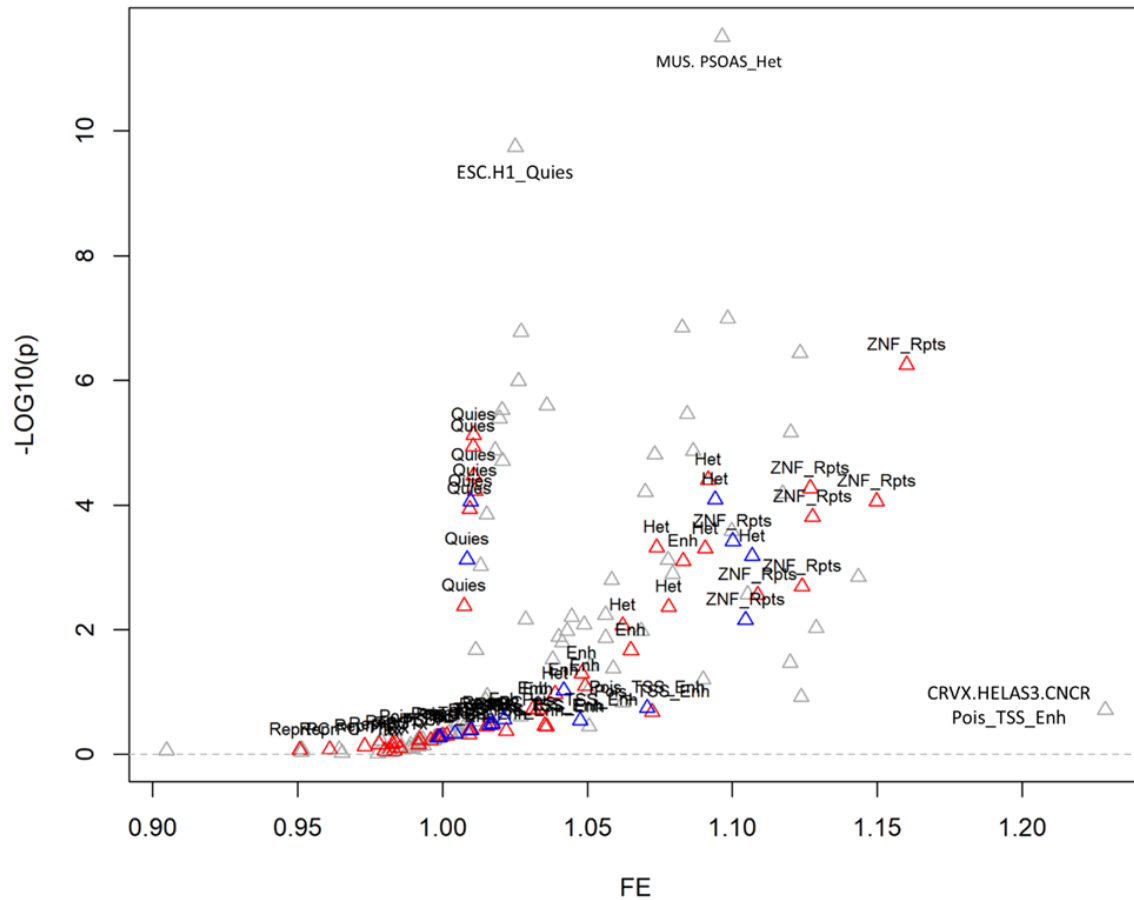


Figure S16 Enrichment analysis of a control SNP set in different chromatin states. Based on eight chromatin states retrieved in two neural crest cell (NCC), six craniofacial tissue (CT) and eleven Roadmap samples, the enrichment of $n=22,999$ SNPs ($P_{MAIC}>0.1$, matched for allele frequency distribution) was calculated using GREGOR (Schmidt et al. 2015). Roadmap samples (in gray) included three fetal (fetal muscle trunk (MUS.TRNK.FET), fetal muscle leg (MUS.LEG.FET), fetal stomach (GI.STMC.FET) and eight non-fetal samples (ESC H1 cell line (ESC.H1), iPSC cell line (IPSC.DF.19.11), bone marrow derived cultured mesenchymal stem cells (STRM.MRW.MSC), primary B cells from peripheral blood (BLD.CD19.PPC), foreskin fibroblast primary cells skin01 (SKIN.PEN.FRSK.FIB.01), psoas muscle (MUS.PSOAS), rectal mucosa donor 29 (GI.RECT.MUC.29), and heLa-S3 cervical carcinoma cell line (CRVX.HELAS3.CNCR)). Enrichment of NCC (blue) and CT (red) samples is indicated by corresponding chromatin states. Roadmap samples with highest $-\log_{10}(p)$ and $\log_2(FE)$ are annotated by tissue type and chromatin state. The chromatin state transcription starting site (TSS) includes both active TSS and upstream/downstream flanking TSS, and the enhancers (Enh) include active and genic enhancers; FE - fold enrichment; ZNF_Rpts - Zink-finger genes and repeats; Het - heterochromatin; Poiss_TSS_Enh - poised enhancers and bivalent TSS. This Figure complements Figure 3b of the Main Text.

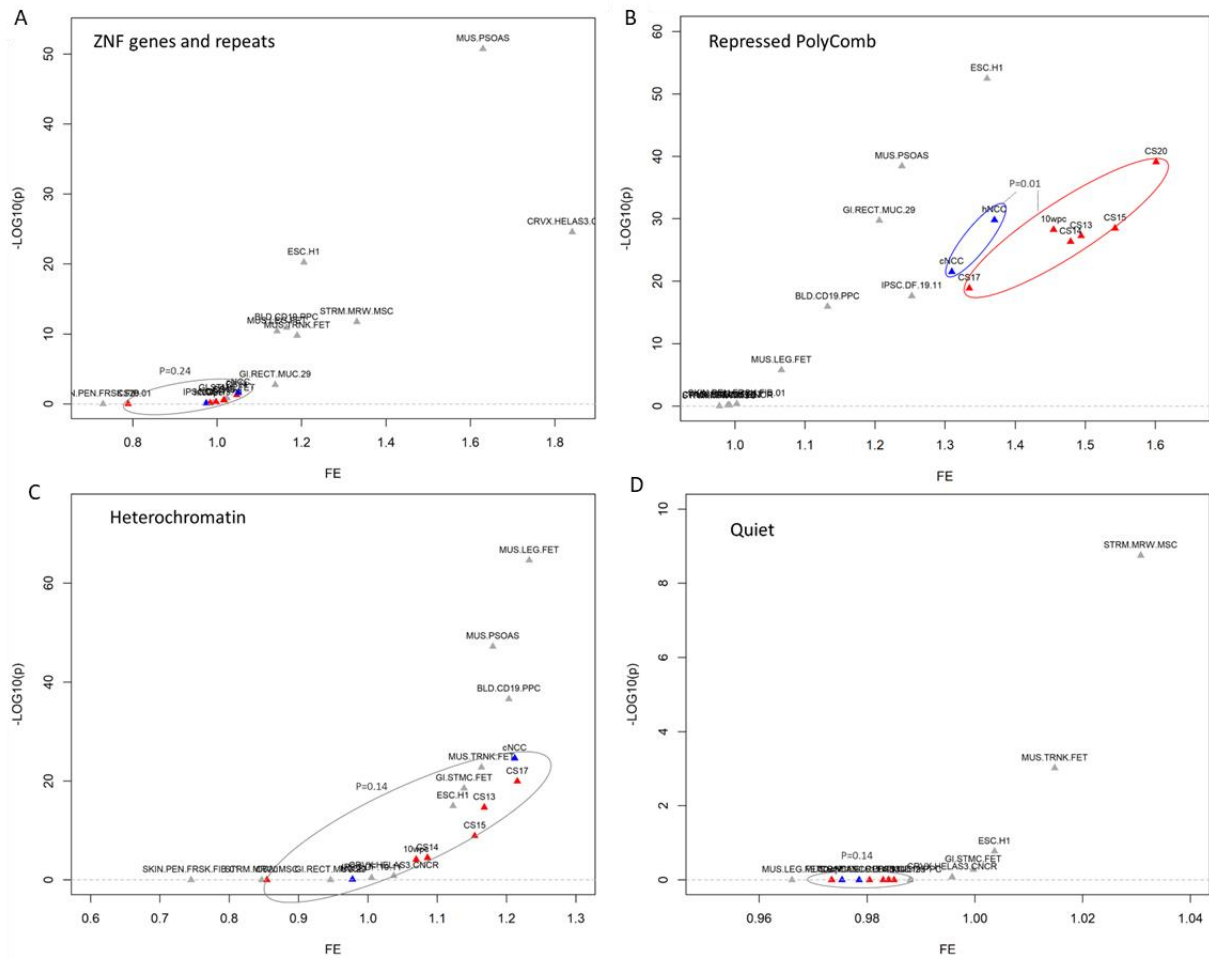


Figure S17 Association of meta-analysis in clefting (MAiC) data for four chromatin states. Based on eight chromatin states retrieved in two neural crest cell (NCC, blue), six craniofacial tissue (CT, pink) and eleven Roadmap samples (gray), enrichment analyses were performed for MAiC risk variants, at $P_{MAiC} \leq 0.001$ ($n=22,999$). Individual enrichment results for MAiC risk variants in four chromatin states. P-values represent difference in enrichment between NCC and CT. **A)** ZNF genes and repeats, **B)** repressed polycomb, **C)** heterochromatin, and **D)** quiescent regions. Abbreviations of tissues as provided by Roadmap. This Figure complements Figures 3c-3f of the Main Text. FE - fold enrichment; ZNF - zinc finger

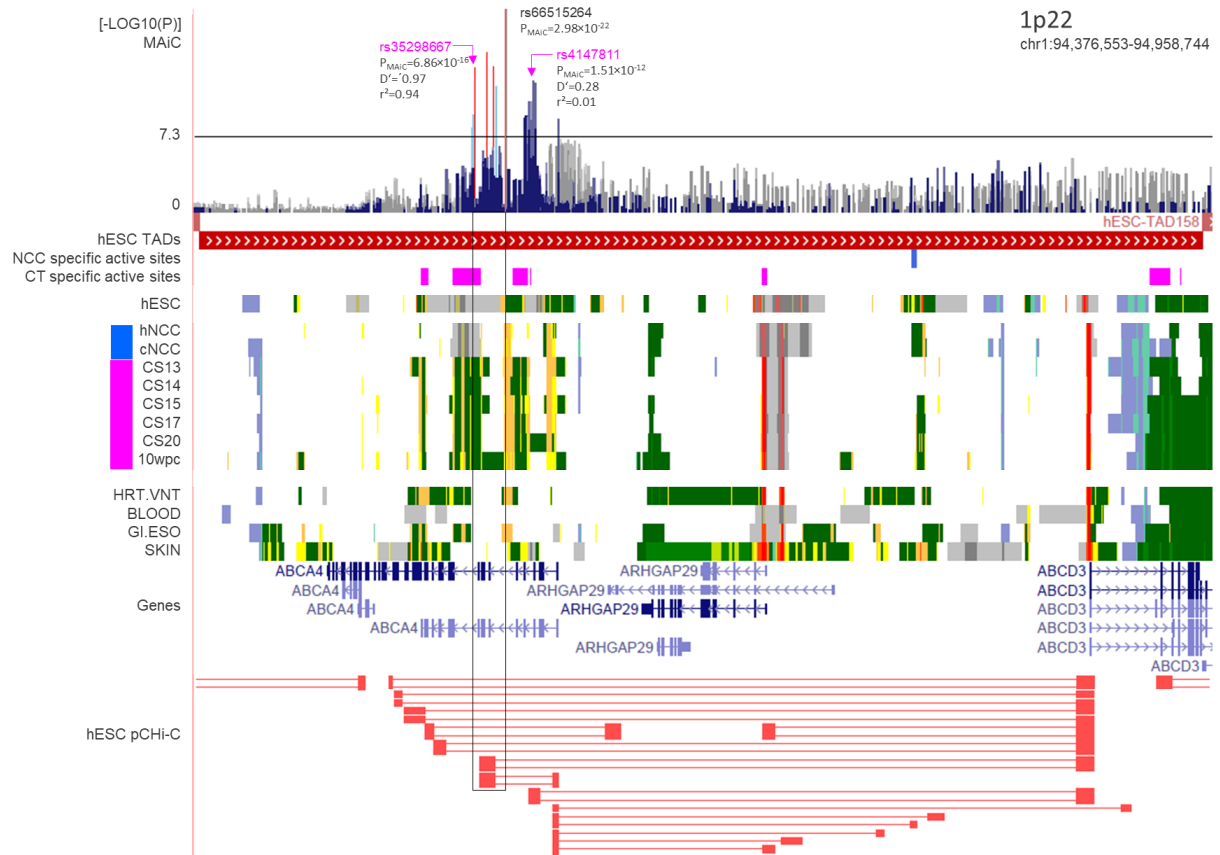


Figure S18 Regulatory architecture at nsCL/P risk locus 1p22. Based on the extent of the topologically associated domain around rs66515264 and variants in linkage disequilibrium (LD) $r^2 \geq 0.6$, different layers of data were aggregated. Tracks include (top-down): MAiC association P-values with colour code based on LD to rs66515264; chromatin segmentation data from early human neural crest cell (hNCC), cranial NCC (cNCC), craniofacial tissue (CT) of different Carnegie stages (CS) and 10 weeks *post-conceptum* (wpc, colour code as in Figure 3 of main text); RefSeq gene positions; and promotor capture (pC) Hi-C cis-interactions collected in hESC (GSE8682133). Black box highlights region of strongest association that contains craniofacial-tissue specific active sites and 3D connections to the promoter of *ABCD3*, which is indicated as active.

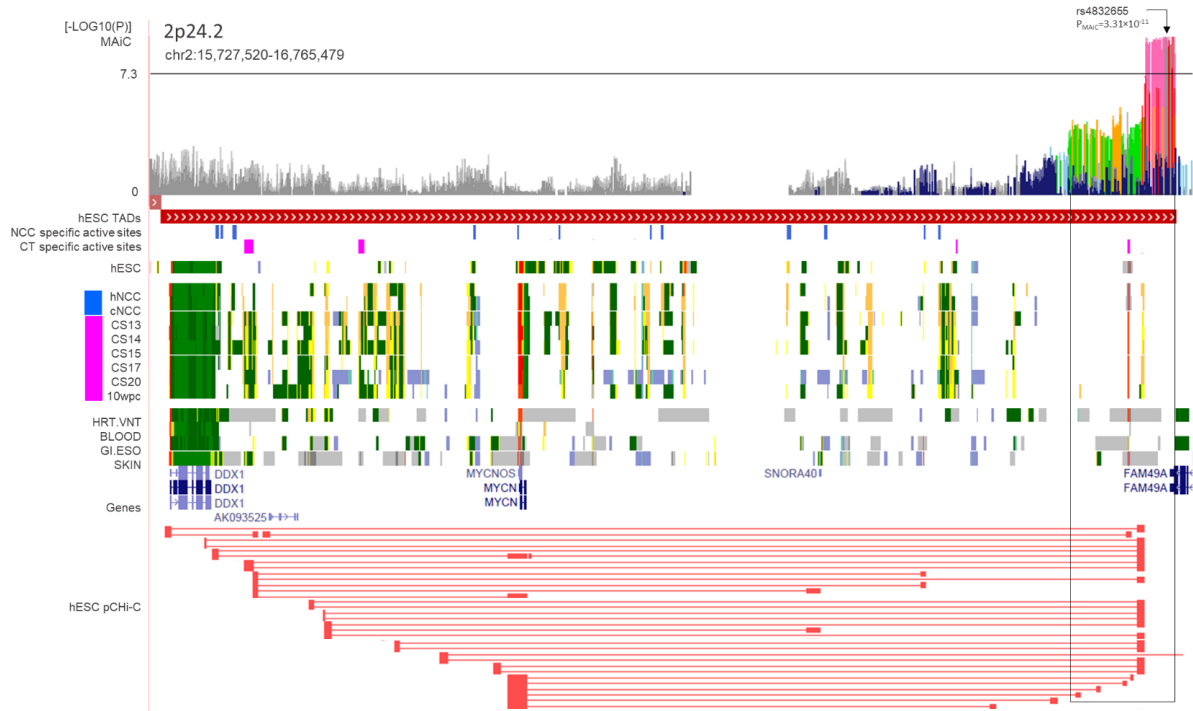


Figure S19 Regulatory architecture at 2p24.2. Based on the extent of the topologically associated domain (TAD) around rs4832655 and variants in linkage disequilibrium (LD) $r^2 \geq 0.6$, different layers of data were aggregated. Tracks include (top-down): MAiC association P-values with colour code based on LD to rs4832655; chromatin segmentation data from early human neural crest cell (hNCC), cranial NCC (cNCC), craniofacial tissue (CT) of different Carnegie stages (CS) and 10 weeks *post-conceptum* (wpc, colour code as in Figure 3 of main text); RefSeq gene positions; and promotor capture Hi-C *cis*-interactions collected in hESC (GSE8682133). Black box highlights region of strongest association that contains SNPs in strong LD to rs4832655 which map to an active region across midfacial development (orange, right) and an active site predominantly in craniofacial tissue (red, left). The presumed enhancer region interacts with diverse genes within the TAD, including *MYCN* and *DDX1*.



Figure S20 Regulatory architecture at 4p13. Based on the extent of the topologically associated domain (TAD) around rs67451576 and variants in linkage disequilibrium (LD) $r^2 \geq 0.6$, different layers of data were aggregated. Tracks include (top-down): MAiC association P-values with colour code based on LD to rs67451576; chromatin segmentation data from early human neural crest cell (hNCC), cranial NCC (cNCC), craniofacial tissue (CT) of different Carnegie stages (CS) and 10 weeks *post-conceptum* (wpc, colour code as in Figure 3 of main text); RefSeq gene positions; and promotor capture (pC) Hi-C cis-interactions collected in hESC (GSE8682133). Black box highlights region of strongest association, which maps to a region of strong craniofacial-tissue specific activity. Chromosomal interactions are indicated to a yet un-characterised genetic region upstream of the *LIMCH1* promoter. Of note is also the lack of expression for the 3'-part of the *LIMCH1* gene in NCC, suggesting the presence of specific isoforms, whose role will have to be further investigated.

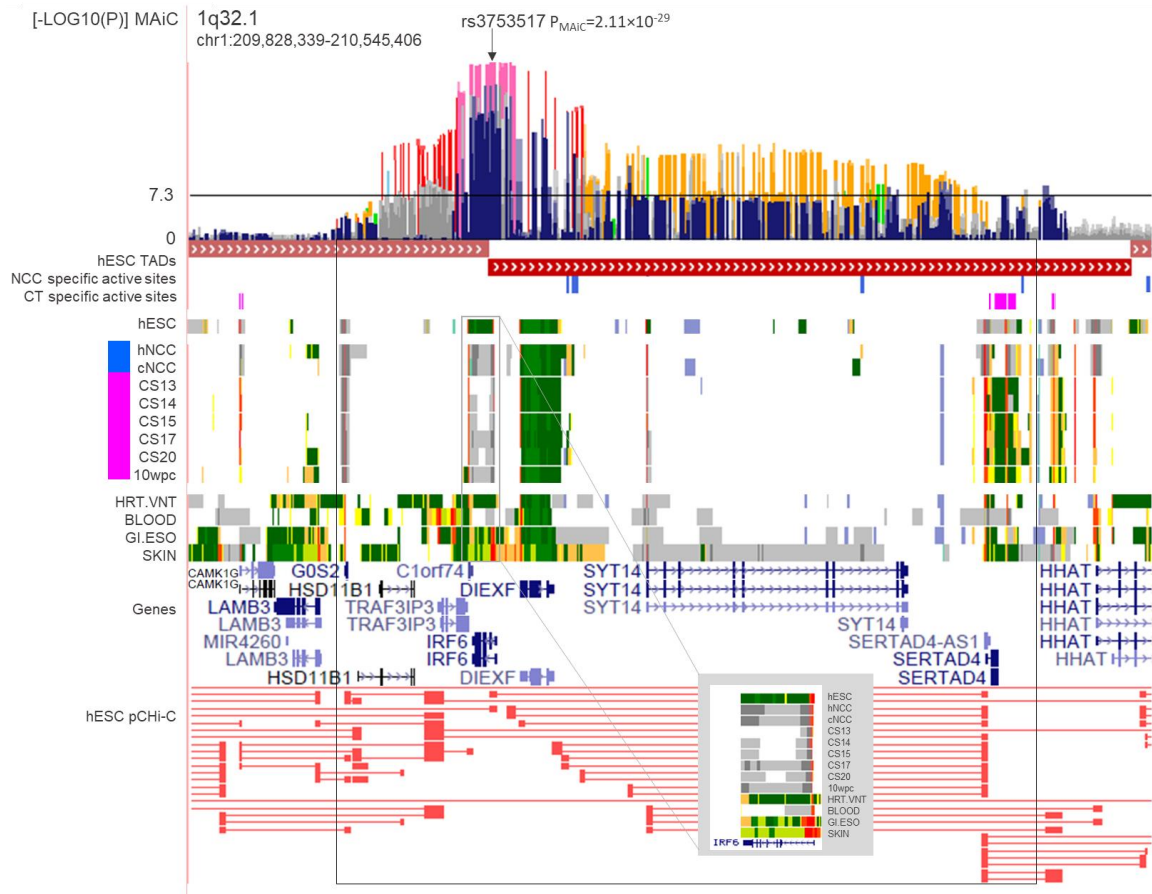


Figure S21 Complex regulatory architecture at 1q32.1. Based on the extent of the topologically associated domain (TAD) around rs3753517 and variants in linkage disequilibrium (LD) $r^2 \geq 0.6$, different layers of data were aggregated. Tracks include (top-down): MAiC association P-values with colour code based on LD to rs3753517; chromatin segmentation data from early human neural crest cell (hNCC), cranial NCC (cNCC), craniofacial tissue (CT) of different Carnegie stages (CS) and 10 weeks *post-conceptum* (wpc, colour code as in Figure 3 of main text); RefSeq gene positions; and promotor capture Hi-C *cis*-interactions collected in hESC (GSE8682133). The region of strongest association maps upstream of the *IRF6* promoter, encompassing a previously identified causal element (Rahimov et al. 2008). Notably, this putative enhancer region is poised in both NCC and CT, which matches the signals of the *IRF6* coding region (grey box). This is likely due to the established function of *IRF6* in periderm / epithelial lineages, which are underrepresented in NCC and CT.

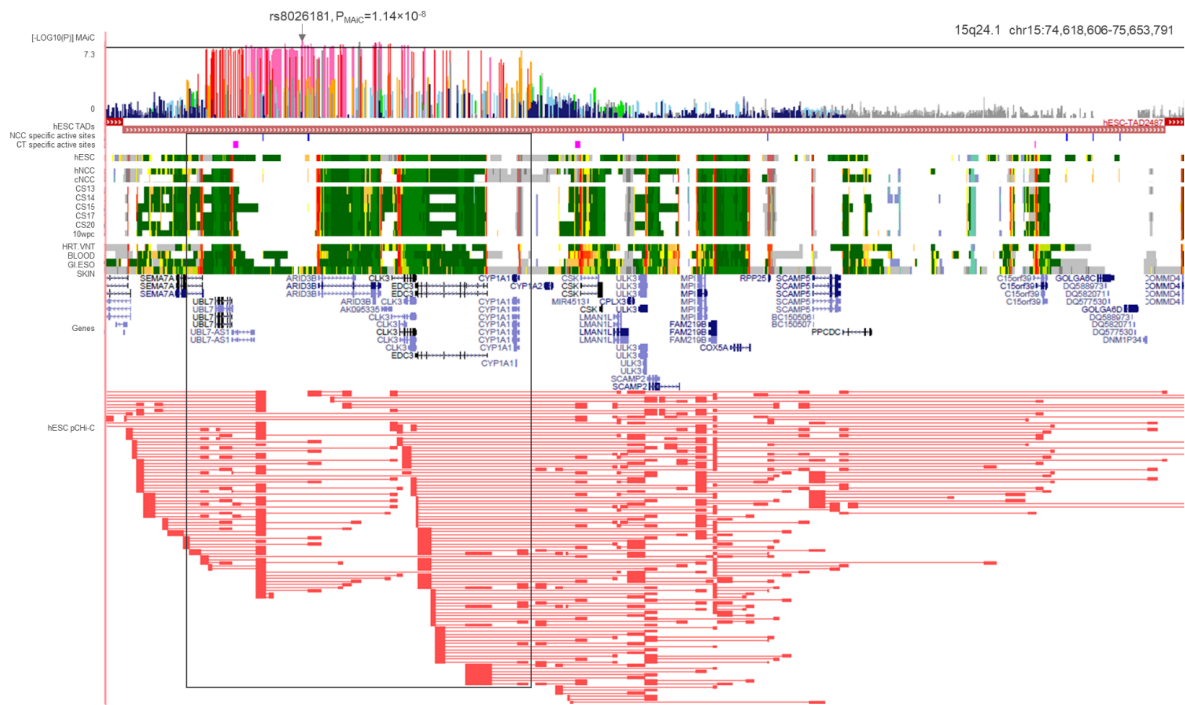


Figure S22 Complex regulatory architecture at 15q24.1. Based on the extent of the topologically associated domain (TAD) around rs8026181 and variants in linkage disequilibrium (LD) $r^2 \geq 0.6$, different layers of data were aggregated. Tracks include (top-down): MAiC association P-values with colour code based on LD to rs8026181; chromatin segmentation data from early human neural crest cell (hNCC), cranial NCC (cNCC), craniofacial tissue (CT) of different Carnegie stages (CS) and 10 weeks *post-conceptum* (wpc, colour code as in Figure 3 of main text); RefSeq gene positions; and promotor capture Hi-C *cis*-interactions collected in hESC (GSE8682133).

Supplemental Tables:

This document contains Tables S1, S3, S5 and S15. All other Supplemental Tables can be found in the Spreadsheet.

Table S1 Overview of individual cohorts used in our meta-analysis in clefting (MAiC).

GWAS (PubMed ID)	Cohort structure	Ethnicity	N individuals (case/control, in trios)	N_{eff} for MAiC (MAiC_{Euro})
Bonn (20023658)	case/control	European	399/1318	1225
GENEVA (20436469)	case-parent trios	European	666	2992 (1332)
		Asian	795	
POFC (27033726)	case/control	European	163/733	2198 (533)
		Asian/Latin American	685/828	
	case-parent trios	European	289	2476 (578)
Asian/Latin American	949			

Three GWAS cohorts comprise four independent sub-cohorts. Abbreviations: N = number after removal of individuals based on inter-individual relationship, N_{eff} = effective number of individuals, calculated as described in "Statistical analysis".

Table S3 Antibodies used in ChIP-Seq experiments and information about imputation.

Cell/tissue ^a	H3K27ac	H3K4me1	H3K4me3	H3K27me3	H3K9me3	H3K36me3	Reference
Prediction	Enhancer	Enhancer	Promotor activating	Repressed	Hetero-chromatin	Active transcription	
hNCC	ab4729, Abcam	ab8895, Abcam	39159, Active Motif	39536, Active Motif	Imputed ¹	Imputed ¹	Rada-Iglesias et al. 2012 (GSE28874)
cNCC	39133, Active Motif						Prescott et al. 2015 (GSE70751)
CT_CS13	ab4729, Abcam	Imputed ²	ab8580, Abcam	07-449, EMD Millipore	Imputed ²	ab9050, Abcam	Wilderman et al. 2018 (GSE97752)
CT_CS14							
CT_CS15							
CT_CS17							
CT_CS20							
10wpc	Imputed ²	Imputed ²	Imputed ²	Imputed ²			

^a - Cell/tissue as origin of chromatin immunoprecipitation followed by sequencing (ChIP-seq); hNCC - early human neural crest cell; cNCC - cranial NCC; CT - craniofacial tissue; CS - Carnegie stage; wpc - weeks *post-conceptum*. ¹ - Imputation performed using ChromImpute v1.0.1 (Ernst and Kellis, 2015) based on 127 cell types from Roadmap Epigenome Project (Roadmap Epigenomics Consortium et al. 2015) and the available chromatin marks in hNCC and cNCC in present study. ² - Imputation performed using ChromImpute v1.0.1 (Ernst and Kellis, 2015) based on 127 cell types from Roadmap Epigenome Project (Roadmap Epigenomics Consortium et al. 2015) and available chromatin marks in CT by Wilderman et al. 2018.

Table S5 Uniquely aligned reads per sample and chromatin mark in neural crest cell samples

Chromatin mark	hNCC^a	cNCC^b
H3K27C	19,673,201	29,813,573
H3K27me3	15,902,493	19,534,085
H3K4me1	17,415,485	24,603,013
H3K4me3	18,482,513	22,506,230
Input	19,733,798	25,999,759
Sum	91,207,490	122,456,660
Average	18,241,498	24,491,332

^a - Number of read in raw data (fastq file) per early human neural crest cell sample (hNCC) and chromatin mark downloaded from GEO (GSE28874); Rada-Iglesias et al. 2012. ^b - Number of read in raw data (fastq file) per cranial NCC (cNCC) sample and chromatin mark downloaded from GEO (GSE70751); Prescott et al. 2015

Table S15 Transfer of 18-state model of ChromHMM to 8-state model.

Original 18 state model				Condensed 8 state model	
States ^a	Description	Color name	RGB code	RGB code	State ^a
TssA	Active TSS	Red	255,0,0	255,0,0	TSS
TssFlnk	Flanking TSS	Orange Red	255,69,0		
TssFlnk	Flanking TSS Upstream	Orange Red	255,69,0		
TssFlnk	Flanking TSS Downstream	Orange Red	255,69,0		
Tx	Strong transcription	Green	0,128,0	0,128,0	Tx
TxWk	Weak transcription	DarkGreen	0,100,0		
EnhG1	Genic enhancer1	GreenYellow	194,225,5	255,255,0	Enh
EnhG2	Genic enhancer2	GreenYellow	194,225,5		
EnhA1	Active Enhancer 1	Orange	255,195,77		
EnhA2	Active Enhancer 2	Orange	255,195,77		
EnhWk	Weak Enhancer	Yellow	255,255,0		
ZNF/Rp	ZNF genes & repeats	Medium Aquamarine	102,205,170	102,205,170	ZNF_Rpts
Het	Heterochromatin	PaleTurquoise	138,145,208	138,145,208	Het
TssBiv	Bivalent/Poised TSS	IndianRed	205,92,92	233,150,122	TssBiv_Enh
EnhBiv	Bivalent Enhancer	DarkKhaki	189,183,107		
ReprPC	Repressed PolyComb	Silver	128,128,128	128,128,128	ReprPC
ReprPC	Weak Repressed PolyComb	Gainsboro	192,192,192		
Quies	Quiescent/Low	White	255,255,255	255,255,255	Quies

a - TSS - transcription starting site; Enh - enhancer; ReprPC - repressed PolyComb; Tx - transcribed sites; Het - Heterochromatin; Pois_TSS_Enh - poised TSS and bivalent enhancers; ZNF_Rpts - Zinc finger genes and repeats.

Supplemental Text - Description of novel risk loci

1) Risk locus 1p36_{CAPZB}

The 1p36 locus was previously suggested as nsCL/P risk locus without reaching formal statistical thresholds, in an association study containing a subsample of the present study³⁰. We here confirm this association at genome-wide significance. The top associated variant in MAiC is rs34746930, located within the genic region of the capping protein (actin filament) muscle Z-line, beta gene (*CAPZB*). Notably, this locus is independent from another risk locus previously

reported at 1p36, around the *PAX7* gene (1p36_{PAX7}, leadSNP rs742071)³¹, as indicated by the location in two different topologically-associated domains (TAD), and the lack of linkage disequilibrium (LD) between the two lead variants (Figure S2); CEU: $r^2=0.0017$ / $D'=0.12$; East Asians: $r^2=0.0058$ / $D'=0.18$; South Asians: $r^2=0.0002$ / $D'=0.16$, assessed using LDpair in LDlink, 1000 genomes phase 3). The TAD around rs34746930 contains several genes, three of which can be considered strong candidates for an involvement in nsCL/P:

CAPZB has been shown to be highly expressed in the first pharyngeal arch during human development, an embryonic structure that hosts cells required for the formation of facial structures³². A *de novo* balanced translocation, disrupting *CAPZB*, was previously reported in a female individual presenting with craniofacial defects (e.g., cleft palate, micrognathia, low-set and rotated ears), hypotonia and developmental delay³³. In addition, several deletions of different sizes encompassing *CAPZB* have been reported in individuals of the DECIPHER database, all of which presented with some degree of craniofacial malformation³⁴. In zebrafish, loss of *capzb* leads to craniofacial phenotypes, and molecular data showed cell migration defects in zebrafish larvae, and differential expression of *pax3a* and *pax7a* in zebrafish neural crest cells³³. Notably, *PAX7* is a candidate gene at the neighbouring 1p36_{PAX7}-locus, suggesting further follow-up-studies including interaction analyses between *CAPZB* and *PAX7*.

NBL1 (neuroblastoma, suppression of tumorigenicity 1) encodes a bone morphogenic protein (BMP) antagonist of the DAN family, the latter is strongly evolutionary conserved. These secreted proteins are involved as antagonists in the BMP pathway: they bind BMP and, thereby, prevent its interaction with other receptors. This suggests an important role during growth and development. Recently, a role for the encoded NBL1-protein in neural crest cell migration was observed: Through *in vivo* analyses in the chicken and *in silico* simulations, it was shown that Nbl1 restrains cell migration through the regulation of cell speed. Thus, NBL1 is suggested to inhibit uncontrolled neural crest invasion and promotes collective migration³⁵.

Finally, the gene **HTR6**, encoding the Serotonin receptor 6, is located ~220kb away from the sentinel SNP. No direct evidence for a role of this gene in craniofacial development has yet been reported. However, it has been suggested that a five-SNP haplotype in *HTR6* is associated with the risk of becoming a smoker³⁶. Given increasing evidence for smoking being an environmental risk factor for orofacial clefting³⁷, this gene might be considered in further analyses of gene-environment analyses (taking into consideration maternal-fetal interactions). However, epigenetic data across mid-facial development do not indicate expression of *HTR6*.

The core associated region (defined as the region containing variants with $r^2>0.8$) extends over 35 kb and contains 25 common variants. One of these variants, rs6682099 (CEU: $r^2=0.92$ / $D'=1.0$ to rs34746930) is highly conserved, has a CADD score >20 and a Regulome-db-Score of 2b. The most prominent position weight matrix (PWM) altered by the C/T exchange of rs6682099 is a 7bp core motif for *PITX2*. Mutations in *PITX2* cause Axenfeld-Rieger Syndrome type 1, an autosomal-dominant disorder affecting primarily facial structures [OMIM 180500]. Patients with Axenfeld-Rieger Syndrome type 1 present with maxillary hypoplasia, short philtrum and thin upper lip, hypodontia (in particular

maxillary incisors) as well as complex eye phenotypes. In GTEx-data (v8), rs6682099 is an eQTL for *NBL1* (in sun-exposed skin tissue, and stomach), *CAPZB* (skin and adrenal gland), and *PQLC2* (in lung tissue), together with a set of other SNPs in high LD.

2) Risk locus 5p12_{FGF10}

The MAiC lead SNP at 5p12, rs60107710, is located about 510 kb away from another previously reported variant at 5p12, i.e. rs10462065³⁸. Although both lead variants are located within the same TAD, there is evidence from haplotype data to be independent from one another, as the lead variants do not share any LD in the three main investigated populations (i.e., CEU: $r^2=0.0044$ / $D'=0.09$; East Asians: $r^2=0.0017$ / $D'=0.08$; South Asians: $r^2=0.0024$ / $D'=0.14$, assessed using LDpair in LDlink, 1000 genomes phase 3 data; Figure S3). The 5p12-associated region around rs10462065 (5p12_{rs10462065}) is located about 320 kb downstream of the *FGF10* transcription start site (TSS), while the newly identified 5p12-association region around rs60107710 (5p12_{rs60107710}) maps about 190 kb upstream of the TSS. The core associated region of the 5p12_{rs60107710}-region contains ~200 SNPs and extends over 180 kb. Within that region, the common variant with the lowest Regulome db-Score is rs1482664 (2b), however, no further annotation is provided in support for this variant being causal (including the absence of a significant eQTL effect in GTex v8).

Within the TAD around rs60107710, two protein-coding genes are located. ***FGF10*** (fibroblast growth factor 10) is a signalling growth factor that predominantly acts in mesenchymal and epithelial tissue, and is required for the development of multiple organs including the craniofacial complex³⁹. *FGF10* is well studied for its role in palatal growth, in particular within the *PAX9* palatogenesis pathway⁴⁰. When *Fgf10* is conditionally knocked-out in murine neural crest cells, many of the phenotypes observed in constitutive *Fgf10*^{-/-} mice are recapitulated, including the frequent occurrence of cleft palate⁴¹. However, in our epigenetic data *FGF10* shows only limited evidence for being actively described in NCC.

The second gene within this TAD, ***NNT*** (nicotinamide nucleotide transhydrogenase), encodes for an integral protein of the inner mitochondrial membrane. It is ubiquitously expressed. Knocking down *NNT* in Hep1-cells results in impaired homeostasis of cells and reduced cell proliferation⁴². So far, no specific role in craniofacial development or disease has been reported.

3) Risk locus 5q13.1

The association at 5q13.1 is characterized by the lead variant rs6449957, which is located ~28 kb upstream of the TSS of ***PIK3R1*** (phosphoinositide-3-Kinase regulatory subunit 1, alias: *GRB1*). However, the associated region extends into the first coding exons of *PIK3R1* (Figure 4c, Figure S4), which is also the only protein-coding gene located within the TAD.

PIK3R1 has been shown to play an important role in the metabolism of insulin^{43,44}. Moreover, heterozygous mutations in the *PIK3R1* gene have been described as causal for SHORT syndrome [OMIM #269880], clinical symptoms of which include teething delay, short stature, hernia and ocular depression⁴⁵. In a recent systems genetics study, *PIK3R1* was identified as candidate gene for nsCL/P based on a re-analysis of a previously published expression dataset that compared dental pulp stem cells from nsCL/P patients with non-affected control children⁴⁶. Moreover, another study employing systems genetics suggested *PIK3R1* as mediator for viral cancerogenesis and interaction with cancer genes⁴⁷, which is noteworthy given some suggestive evidence of orofacial clefting being associated with an increased risk of different cancer types⁴⁸.

The association structure at 5q13.1 is described in the Main Text. Briefly, the core associated region comprises ~30 variants, none of which is an eQTL in GTEx v8. However, we observed rs6449957 to be reported as splice QTL for both *PIK3R1* and a long non-coding RNA (*LINC02219*), respectively, in testis. The lowest Regulome-db score at 5p13.1 is observed for rs6449956 (score 2b), which is predicted to disrupt the binding site for transcription factor FEV, a member of the Ets-family of transcription factors (according to JASPAR2018, Figure S4), however, no role of FEV in craniofacial development has yet been reported. Interestingly, our integrative data suggests some evidence for *MAST4*, located in the adjacent TAD, as second candidate gene at this locus (see Main Text).

4) Risk locus 7p21.1

The 7p21.1 risk region is characterized by its lead SNP rs62453366, which is located intronically within the *ABCB5* gene (Figure S5). The core associated region is very narrow, encompassing 10 kb only. None of the 20 variants located within the GWAS_{SNP}-region is an eQTL in GTEx v8 data, and none of the variants has CADD>10 or Regulome-score better than 4. Within the TAD, three genes are located – *ABCB5*, *SP8* and *RPLS23P8*, the latter of which is a processed ribosomal protein pseudogene for which no functional information is available. We therefore here describe the two other genes, which both can be considered interesting candidate genes for nsCL/P.

ABCB5 (ATP binding cassette subfamily B member 5) represents a member of the ABC transporter superfamily of integral membrane proteins. *ABCB5* is a marker for progenitor cells in both skin and human melanoma, and plays an important role as regulator of cellular differentiation⁴⁹. *ABCB5* is also expressed in specific fractions of limbal stem (LS) cells, lack of which represent a major cause of blindness⁵⁰. LS-cells positive for *ABCB5* expression co-express the *deltaNp63alpha isoform* of p63, but not the differentiation marker *KRT12*⁵¹. The striking co-expression of *ABCB5* and p63 in limbal stem cells is noteworthy, for the following reasons: (i) p63 is the causal gene for different types of ectodermal dysplasias (e.g. AEC/EEC-syndrome), (ii) the *p63* gene maps to an nsCL/P risk locus itself (chromosome 3q28), and (ii) p63 has recently been shown to be critically relevant for establishing enhancer marks at genes relevant in craniofacial development and disease⁵³. Notably, a recent study identified that p63 is superior to *ABCB5* as marker for stem cells, and also associates with LS-cells with increased pigmentation⁵³.

SP8 encodes the zinc finger transcription factor SP8. *Sp8*^{-/-}-mice show severe defects in limb development and craniofacial malformations. Specifically, at E14.5, facial prominences of *Sp8*^{-/-} mice are underdeveloped in both size and structure, which results in severe craniofacial hypoplasia. Later in development, this is still recapitulated through severe midline defects, exencephaly, cleft palate, and a loss of neural crest cell. In that study, *Sp8* was identified as craniofacial signalling center that regulates proliferation and apoptosis of NCC, with molecular downstream effects on *Fgf8* and *Fgf17* expression. Partial rescue of the phenotype in the *Sp8*^{-/-} mice was obtained through reduction of Sonic hedgehog signalling, indicative of role for *Sp8* in the Shh-Fgf signalling pathway⁵⁴. Recently, a truncating mutation within *SP8* was identified in a patient with nsCL/P (p.S261X) in a resequencing study⁵⁵, which is highly interesting as no *SP8* loss-of-function variant is currently reported in gnomAD database.

5) Risk locus 20q13.12

The 20q13.12 risk locus is characterized by its lead SNP rs3091552, which is located upstream of the gene eyes absent homologue 2 (***EYA2***). *EYA2* is also the only protein-coding gene located in this single-gene TAD (Figure S6). The core-associated region encompasses rs3091552 and two additional variants

in strong LD (rs12481092: $D' = 1.0 / r^2 = 0.88$; rs6066089: $D' = 1.0 / r^2 = 0.83$). In GTEx (v8), all three SNPs represent eQTLs for *EYA2* in artery/aorta and additional tissues, with the risk alleles being associated with decreased *EYA2* expression. The strongest eQTL effect at this locus is demonstrated for rs8125695 ($D' = 1.0 / r^2 = 0.79$ to rs3091552, with $P=8.6 \times 10^{-14}$, effect size 0.48 in GTExv8).

The gene *EYA2* encodes for a transcription factor with profound role in a variety of cellular and developmental processes, including cardiac⁵⁶ and muscle⁵⁷ development. Members of the EYA-family (including *EYA2*) are centrally involved in embryonic organogenesis through the promotion of proliferation and/or survival of progenitor-cell populations⁵⁸, and loss of function mutations have been shown to cause branchio-oto-renal (BOR) syndrome⁵⁹ which, among others, is characterized by malformations of anatomical structures derived from the human branchial arches (e.g., ears, OMIM: #113650). In mice it was shown that during eye morphogenesis, retinoic acid targets the neural crest-cell-derived mesenchyme in which *Eya2*-related apoptosis has been observed⁶⁰. EYA-proteins largely functions through formation of a protein-complex together with SIX1 and DACH⁵⁸, and the important role for this transcription complex has been shown through both functional assays⁵⁸ and structural modelling⁶¹, respectively. Notably, disruption of this process contributes to epithelial-mesenchymal-transition, and metastasis⁶¹. Recently, it was also shown in mice that disrupting *Eya2* phosphatase activity through chemical inhibits *Eya2*-mediated cell migration⁶².

References

1. Mangold, E., Ludwig, K.U., Birnbaum, S., Baluardo, C., Ferrian, M., Herms, S., Reutter, H., De Assis, N.A., Chawa, T. Al, Mattheisen, M., et al. (2010). Genome-wide association study identifies two susceptibility loci for nonsyndromic cleft lip with or without cleft palate. *Nat. Genet.*
2. Beaty, T.H., Murray, J.C., Marazita, M.L., Munger, R.G., Ruczinski, I., Hetmanski, J.B., Liang, K.Y., Wu, T., Murray, T., Fallin, M.D., et al. (2010). A genome-wide association study of cleft lip with and without cleft palate identifies risk variants near MAFB and ABCA4. *Nat. Genet.*
3. Leslie, E.J., Carlson, J.C., Shaffer, J.R., Feingold, E., Wehby, G., Laurie, C.A., Jain, D., Laurie, C.C., Doheny, K.F., McHenry, T., et al. (2016). A multi-ethnic genome-wide association study identifies novel loci for non-syndromic cleft lip with or without cleft palate on 2p 24.2, 17q23 and 19q13. *Hum. Mol. Genet.*
4. Ludwig, K.U., Mangold, E., Herms, S., Nowak, S., Reutter, H., Paul, A., Becker, J., Herberz, R., Alchawa, T., Nasser, E., et al. (2012). Genome-wide meta-analyses of nonsyndromic cleft lip with or without cleft palate identify six new risk loci. *Nat. Genet.*
5. Ludwig, K.U., Böhmer, A.C., Bowes, J., Nikolić, M., Ishorst, N., Wyatt, N., Hammond, N.L., Gözl, L., Thieme, F., Barth, S., et al. (2017). Imputation of orofacial clefting data identifies novel risk loci and sheds light on the genetic background of cleft lip ± cleft palate and cleft palate only. *Hum. Mol. Genet.*
6. Leslie, E.J., Carlson, J.C., Shaffer, J.R., Butali, A., Buxó, C.J., Castilla, E.E., Christensen, K., Deleyiannis, F.W.B., Leigh Field, L., Hecht, J.T., et al. (2017). Genome-wide meta-analyses of nonsyndromic orofacial clefts identify novel associations between FOXE1 and all orofacial clefts, and TP63 and cleft lip with or without cleft palate. *Hum. Genet.*
7. Manichaikul, A., Mychaleckyj, J.C., Rich, S.S., Daly, K., Sale, M., and Chen, W.M. (2010). Robust relationship inference in genome-wide association studies. *Bioinformatics.*

8. Marchini, J., Howie, B., Myers, S., McVean, G., and Donnelly, P. (2007). A new multipoint method for genome-wide association studies by imputation of genotypes. *Nat. Genet.*
9. Spielman, R.S., McGinnis, R.E., and Ewens, W.J. (1993). Transmission test for linkage disequilibrium: the insulin gene region and insulin-dependent diabetes mellitus (IDDM). *Am. J. Hum. Genet.*
10. Willer, C.J., Li, Y., and Abecasis, G.R. (2010). METAL: Fast and efficient meta-analysis of genomewide association scans. *Bioinformatics.*
11. Bulik-Sullivan, B., Loh, P.R., Finucane, H.K., Ripke, S., Yang, J., Patterson, N., Daly, M.J., Price, A.L., Neale, B.M., Corvin, A., et al. (2015). LD score regression distinguishes confounding from polygenicity in genome-wide association studies. *Nat. Genet.*
12. Watanabe, K., Taskesen, E., Van Bochoven, A., and Posthuma, D. (2017). Functional mapping and annotation of genetic associations with FUMA. *Nat. Commun.*
13. Machiela, M.J., and Chanock, S.J. (2015). LDlink: A web-based application for exploring population-specific haplotype structure and linking correlated alleles of possible functional variants. *Bioinformatics.*
14. Aguet, F., Brown, A.A., Castel, S.E., Davis, J.R., He, Y., Jo, B., Mohammadi, P., Park, Y.S., Parsana, P., Segrè, A. V., et al. (2017). Genetic effects on gene expression across human tissues. *Nature.*
15. de Leeuw, C.A., Mooij, J.M., Heskes, T., and Posthuma, D. (2015). MAGMA: Generalized Gene-Set Analysis of GWAS Data. *PLoS Comput. Biol.*
16. Rada-Iglesias, A., Bajpai, R., Prescott, S., Brugmann, S.A., Swigut, T., and Wysocka, J. (2012). Epigenomic annotation of enhancers predicts transcriptional regulators of human neural crest. *Cell Stem Cell.*
17. Prescott, S.L., Srinivasan, R., Marchetto, M.C., Grishina, I., Narvaiza, I., Selleri, L., Gage, F.H., Swigut, T., and Wysocka, J. (2015). Enhancer Divergence and cis-Regulatory Evolution in the Human and Chimpanzee Neural Crest. *Cell.*
18. Wilderman, A., VanOudenhove, J., Kron, J., Noonan, J.P., and Cotney, J. (2018). High-Resolution Epigenomic Atlas of Human Embryonic Craniofacial Development. *Cell Rep.*
19. Bajpai, R., Chen, D.A., Rada-Iglesias, A., Zhang, J., Xiong, Y., Helms, J., Chang, C.P., Zhao, Y., Swigut, T., and Wysocka, J. (2010). CHD7 cooperates with PBAF to control multipotent neural crest formation. *Nature.*
20. Rada-Iglesias, A., Bajpai, R., Swigut, T., Brugmann, S.A., Flynn, R.A., and Wysocka, J. (2011). A unique chromatin signature uncovers early developmental enhancers in humans. *Nature.*
21. Furlan-Magaril, M., Rincón-Arango, H., and Recillas-Targa, F. (2009). Sequential chromatin immunoprecipitation protocol: ChIP-reChIP. *Methods Mol. Biol.*
22. Langmead, B., and Salzberg, S.L. (2012). Fast gapped-read alignment with Bowtie 2. *Nat. Methods.*
23. Landt, S.G., Marinov, G.K., Kundaje, A., Kheradpour, P., Pauli, F., Batzoglou, S., Bernstein, B.E., Bickel, P., Brown, J.B., Cayting, P., et al. (2012). ChIP-seq guidelines and practices of the ENCODE and modENCODE consortia. *Genome Res.*
24. Zhang, Y., Liu, T., Meyer, C.A., Eeckhoute, J., Johnson, D.S., Bernstein, B.E., Nussbaum, C., Myers, R.M., Brown, M., Li, W., et al. (2008). Model-based analysis of ChIP-Seq (MACS). *Genome Biol.*

25. Ernst, J., and Kellis, M. (2015). Large-scale imputation of epigenomic datasets for systematic annotation of diverse human tissues. *Nat. Biotechnol.*
26. Roadmap Epigenomics Consortium, Kundaje, A., Meuleman, W., Ernst, J., Bilenky, M., Yen, A., Heravi-Moussavi, A., Kheradpour, P., Zhang, Z., Wang, J., et al. (2015). Integrative analysis of 111 reference human epigenomes. *Nature.*
27. Ernst, J., and Kellis, M. (2017). Chromatin-state discovery and genome annotation with ChromHMM. *Nat. Protoc.*
28. Laugsch, M., Bartusel, M., Alirzayeva, H., Karaolidou, A., Rehimi, R., Crispatzu, G., Nikolic, M., Bleckwehl, T., Kolovos, P., van Ijcken, W.F.J., et al. (2018). Disruption of the TFAP2A Regulatory Domain Causes Banchio-Oculo-Facial Syndrome (BOFS) and Illuminates Pathomechanisms for Other Human Neurocristopathies.
29. Schmidt, E.M., Zhang, J., Zhou, W., Chen, J., Mohlke, K.L., Chen, Y.E., and Willer, C.J. (2015). GREGOR: Evaluating global enrichment of trait-associated variants in epigenomic features using a systematic, data-driven approach. *Bioinformatics.*
30. Carlson, J. C. *et al.* (2019). A systematic genetic analysis and visualization of phenotypic heterogeneity among orofacial cleft GWAS signals. *Genet Epidemiol* 43, 704-716.
31. Ludwig, K. U. *et al.* (2012). Genome-wide meta-analyses of nonsyndromic cleft lip with or without cleft palate identify six new risk loci. *Nat Genet* 44, 968-971.
32. Cai, J. *et al.* (2005). Gene expression in pharyngeal arch 1 during human embryonic development. *Hum Mol Genet* 14, 903-912.
33. Mukherjee, K. *et al.* (2016). Actin capping protein CAPZB regulates cell morphology, differentiation, and neural crest migration in craniofacial morphogenesis. *Hum Mol Genet* 25, 1255-1270.
34. Kang, S. H. *et al.* (2007). Identification of proximal 1p36 deletions using array-CGH: a possible new syndrome. *Clin Genet* 72, 329-338.
35. McLennan, R. *et al.* (2017). DAN (NBL1) promotes collective neural crest migration by restraining uncontrolled invasion. *J Cell Biol* 216, 3339-3354.
36. Lerer, E., Kanyas, K., Karni, O., Ebstein, R. P. & Lerer, B. (2006). Why do young women smoke? II. Role of traumatic life experience, psychological characteristics and serotonergic genes. *Mol Psychiatry* 11, 771-781.
37. Sabbagh, H. J. *et al.* (2015). Passive smoking in the etiology of non-syndromic orofacial clefts: a systematic review and meta-analysis. *PLoS One* 10, e0116963.
38. Yu, Y. *et al.* (2017). Genome-wide analyses of non-syndromic cleft lip with palate identify 14 novel loci and genetic heterogeneity. *Nat Commun* 8, 14364.
39. Prochazkova, M., Prochazka, J., Marangoni, P. & Klein, O. D. (2018). Bones, Glands, Ears and More: The Multiple Roles of FGF10 in Craniofacial Development. *Front Genet* 9, 542.
40. Li, R., Chen, Z., Yu, Q., Weng, M. & Chen, Z. (2019). The Function and Regulatory Network of Pax9 Gene in Palate Development. *J Dent Res* 98, 277-287.
41. Teshima, T. H., Lourenco, S. V. & Tucker, A. S. (2016). Multiple Cranial Organ Defects after Conditionally Knocking Out Fgf10 in the Neural Crest. *Front Physiol* 7, 488.

42. Ho, H. Y., Lin, Y. T., Lin, G., Wu, P. R. & Cheng, M. L. (2017). Nicotinamide nucleotide transhydrogenase (NNT) deficiency dysregulates mitochondrial retrograde signaling and impedes proliferation. *Redox Biol* 12, 916-928.
43. Thauvin-Robinet, C. *et al.* (2013). PIK3R1 mutations cause syndromic insulin resistance with lipoatrophy. *Am J Hum Genet* 93, 141-149.
44. Kuo, T. *et al.* (2017). Pik3r1 Is Required for Glucocorticoid-Induced Perilipin 1 Phosphorylation in Lipid Droplet for Adipocyte Lipolysis. *Diabetes* 66, 1601-1610.
45. Avila, M. *et al.* (2016). Clinical reappraisal of SHORT syndrome with PIK3R1 mutations: toward recommendation for molecular testing and management. *Clin Genet* 89, 501-506.
46. Kobayashi, G. S. *et al.* (2013). Susceptibility to DNA damage as a molecular mechanism for non-syndromic cleft lip and palate. *PLoS One* 8, e65677.
47. Wang, H. *et al.* (2016). Gene expression profiling analysis contributes to understanding the association between non-syndromic cleft lip and palate, and cancer. *Mol Med Rep* 13, 2110-2116.
48. Bille, C. *et al.* (2005). Cancer risk in persons with oral cleft--a population-based study of 8,093 cases. *Am J Epidemiol* 161, 1047-1055.
49. Frank, N. Y. *et al.* (2003). Regulation of progenitor cell fusion by ABCB5 P-glycoprotein, a novel human ATP-binding cassette transporter. *J Biol Chem* 278, 47156-47165.
50. Dua, H. S., Joseph, A., Shanmuganathan, V. A. & Jones, R. E. (2003). Stem cell differentiation and the effects of deficiency. *Eye (Lond)* 17, 877-885.
51. Ksander, B. R. *et al.* (2014). ABCB5 is a limbal stem cell gene required for corneal development and repair. *Nature* 511, 353-357.
52. Lin-Shiao, E. *et al.* (2019). p63 establishes epithelial enhancers at critical craniofacial development genes. *Sci Adv* 5, eaaw0946.
53. Liu, L. *et al.* (2018). Pigmentation Is Associated with Stemness Hierarchy of Progenitor Cells Within Cultured Limbal Epithelial Cells. *Stem Cells* 36, 1411-1420.
54. Kasberg, A. D., Brunskill, E. W. & Steven Potter, S. (2013). SP8 regulates signaling centers during craniofacial development. *Dev Biol* 381, 312-323.
55. Marini, N. J., Asrani, K., Yang, W., Rine, J. & Shaw, G. M. (2019). Accumulation of rare coding variants in genes implicated in risk of human cleft lip with or without cleft palate. *Am J Med Genet A* 179, 1260-1269.
56. Lee, S. H. *et al.* (2009). The transcription factor Eya2 prevents pressure overload-induced adverse cardiac remodeling. *J Mol Cell Cardiol* 46, 596-605.
57. Grifone, R. *et al.* (2007). Eya1 and Eya2 proteins are required for hypaxial somitic myogenesis in the mouse embryo. *Dev Biol* 302, 602-616.
58. Li, X. *et al.* (2003). Eya protein phosphatase activity regulates Six1-Dach-Eya transcriptional effects in mammalian organogenesis. *Nature* 426, 247-254.
59. Abdelhak, S. *et al.* (1997). A human homologue of the Drosophila eyes absent gene underlies branchio-oto-renal (BOR) syndrome and identifies a novel gene family. *Nat Genet* 15, 157-164.

60. Matt, N. *et al.* (2005). Retinoic acid-dependent eye morphogenesis is orchestrated by neural crest cells. *Development* 132, 4789-4800.
61. Patrick, A. N. *et al.* (2013). Structure-function analyses of the human SIX1-EYA2 complex reveal insights into metastasis and BOR syndrome. *Nat Struct Mol Biol* 20, 447-453.
62. Krueger, A. B. *et al.* (2014). Allosteric inhibitors of the Eya2 phosphatase are selective and inhibit Eya2-mediated cell migration. *J Biol Chem* 289, 16349-16361.