

HDMAX2: A framework for High Dimensional Mediation Analysis with application to maternal smoking, placental DNA methylation and birth outcomes

B. Jumentier (1), C-C. Barrot (1), M. Estavoyer (1), J. Tost (2), B. Heude (3), O. François (1,4,*), J. Lepeule (5,*)

Author affiliations

(1) Université Grenoble-Alpes, Centre National de la Recherche Scientifique, Grenoble INP, TIMC CNRS UMR 5525, 38000 Grenoble, France.

(2) Laboratory for Epigenetics and Environment, Centre National de Recherche en Genomique Humaine, CEA – Institut de Biologie François Jacob, University Paris Saclay, Evry, France

(3) Université de Paris, CRESS, Inserm, INRAE, F-75004 Paris, France.

(4) Inria Grenoble - Rhône-Alpes Inovallée 655 Avenue de l'Europe - CS 90051 38334 Montbonnot, France.

(5) Université Grenoble-Alpes, INSERM, CNRS, Team of Environmental Epidemiology Applied to Development and Respiratory Health, Institute for Advanced Biosciences, 38000 Grenoble, France

(*) These authors contributed equally to the study

Corresponding authors:

Olivier.Francois@univ-grenoble-alpes.fr, Johanna.Lepeule@univ-grenoble-alpes.

Abstract

Most high dimensional mediation epigenetic studies evaluate each mediator indirect effects independently. These approaches are underpowered to detect multiple mediators and the overall indirect effect remains poorly quantified. We developed HDMAX2, an efficient algorithm for high dimensional mediation analysis using max-squared tests, considering CpGs and aggregated mediator regions. HDMAX2 outperformed existing approaches in simulated and real data, detected true mediators with increased power, and assessed the overall indirect effect of several mediators from a high dimension matrix. We showed a polygenic architecture for placenta CpGs and regions mediating the effects of maternal smoking during pregnancy on baby's birth weight.

1. Background

Mediation analysis is a statistical tool used to gain insights into the causal mechanisms that relate an exposure to an outcome (Baron and Kenny, 1986). It is increasingly used in environmental epidemiology, in particular in Developmental Origins of Health and Disease (DOHaD) research and in molecular epidemiology studies (Blum et al., 2020; Nakamura et al., 2021). With the development of high-throughput screening technologies, these methods have become key tools to investigate the pathways by which environmental exposures can affect health outcomes, and more specifically those involving epigenetic mechanisms such as DNA methylation (DNAm) variations (MacKinnon et al., 2007; VanderWeele, 2015, 2016; Zeng et al., 2021). High-dimensional mediation analysis is an extension of unidimensional mediation analysis including multiple mediators (Blum et al., 2020).

A typical high-dimensional analysis for DNAm markers generally includes three main steps. The first step tests both the effects of exposure on DNAm levels and the effects of DNAm levels on the health outcome based on epigenome-wide association studies (EWAS). The second step combines significance values obtained from the two EWAS at the first step in order to perform mediation tests, and assesses the mediator status of each marker. The third step quantify the indirect effects of exposure on the health outcome through DNAm changes. Analyses involving a high dimensional set of mediators are difficult, and raise numerous statistical issues (Blum et al., 2020). To overcome those issues, several approaches have been proposed during the recent years. Classic approaches perform multi-

dimensional analysis by running unidimensional mediation analyses for each DNAm marker, for example using Sobel tests or by estimating Average Causal Mediated Effects (ACME) (Imai et al., 2010; Sobel, 1982). Improvements of the Sobel test for indirect effects combine the significance values obtained from the two EWAS in various ways (Dai et al., 2020; Djordjilović et al., 2020, 2019; Gao et al., 2019; Sampson et al., 2018; Zhang et al., 2016). However, there is no consensus on the most relevant combination of EWAS and mediation tests for a high-dimensional analysis. Furthermore, the overall indirect effect of multiple mediators remains poorly quantified from estimates of single mediator effects in a context of correlation among the mediators.

We addressed the above issues by developing HDMAX2, a method for high-dimensional mediation analysis, and systematically compared HDMAX2 to recently proposed approaches. HDMAX2 relies on latent factor regression models to evaluate associations of exposure and outcome with DNAm, and on mediation tests that control the type I error when combining the significance values obtained in the exposure and outcome EWAS. We developed additional features to further consider differential methylation regions as mediators and to estimate an overall mediated effect of DNAm accounting simultaneously for all mediators identified. Using simulations, we performed an in-depth evaluation of the statistical performances of HDMAX2 in comparison to recently developed methods, and showed that HDMAX2 has increased power compared to the other methods.

We then used HDMAX2 to identify the causal role of placental DNA methylation in the pathway between maternal smoking during pregnancy, gestational age at delivery and birth weight of the baby. While most of the literature has focused on cord blood (Agha et al., 2016; Joubert et al., 2012; Küpers et al., 2015; Xu et al., 2021), few studies have investigated placental DNAm, although it plays a key role in fetal programming (Cardenas et al., 2019; Morales et al., 2016; Rousseaux et al., 2019; Thornburg and Marshall, 2015). We discovered new CpG mediators and regions that were not identified by classic approaches (Cardenas et al., 2019; Morales et al., 2016), and estimated the overall indirect effect of maternal smoking during pregnancy on newborn birth weight and on gestational age at delivery. The reasoning we used for placental DNAm data also holds for other types of tissue and quantitative omics data, and HDMAX2 extends to various data.

2. Methods

2.1 Overview of the HDMAX2 method

HDMAX2 is a newly developed approach for high dimensional mediation analysis structured in three main steps (Figure 1). The first step of HDMAX2 corresponds to an extension of regression models considered generally in unidimensional mediation analysis (Baron and Kenny, 1986). The extension includes latent factors as covariates in the models to account for unobserved variables that confound multidimensional DNAm data analysis, such as batch effects or cell-type heterogeneity in samples. The second step identifies potential mediators by combining paired significance values that are obtained when testing the effect of exposure on DNAm and the effect of DNAm on outcome in step 1. Step 2 is not restricted to CpG markers, and can also identify aggregated mediator regions (AMRs) based on the paired P -values. The third step quantifies indirect effects of identified mediators either separately or simultaneously with an overall indirect effect.

Figure 1: Workflow of the high dimensional mediation analysis procedure implemented in HDMAX2.

Step 1. Evaluating associations between exposure, mediators and outcome. The first step of HDMAX2 consists in adjusting latent factor mixed models (LFMMs) to estimate the effects of exposure, X , on a matrix M of CpG markers, and the effect of each marker on outcome, Y (Caye et al., 2019; Jumentier et al., 2020). LFMMs belong to a class of estimation algorithms that adjust latent factor models, and that encompasses surrogate variable analysis (SVA) (Leek and Storey, 2007), directed SVA (Lee et al., 2017) or confounder adjusted testing and estimation (CATE) (Wang et al., 2017). Latent factor models differ from models based on a priori estimates of cell types (Houseman et al., 2016; Rahmani et al., 2016), and represent a more general approach to the issue of confounding in association studies (Leek and Storey, 2007). Within the latent factor regression framework, additional known covariates (like maternal age or sex of the newborn) can be included in the model to improve accuracy.

Mediation analysis based on DNAm markers performs two EWAS. To estimate the effects of exposure (X) on a matrix of CpG markers (M), the following model was first adjusted to the centered data

$$\mathbf{M} = \mathbf{X}\mathbf{a}^T + \mathbf{U}_1\mathbf{V}_1^T + \mathbf{E}_1, \quad (\text{Equation 1})$$

where \mathbf{a} contains the vector of effect sizes of exposure on DNAm levels, \mathbf{U}_1 is a matrix formed of K latent factors estimated simultaneously with \mathbf{a} , \mathbf{V}_1 contains the loadings associated with the latent factors, and \mathbf{E}_1 is a matrix of residual errors. The K latent factors represent hidden confounders, for example unobserved cell types of tissue samples and batch effects. Using the latent factor regression defined in Eq.1, a significance value, P_x , is computed for the test of a null effect size for exposure on DNAm at each CpG marker ($H_0: a_j = 0$, for the j^{th} marker).

A second EWAS was then performed in order to estimate effect sizes for the DNAm levels on the health outcome (Y) as follows:

$$\mathbf{Y} = \mathbf{Xc} + \mathbf{Mb}^T + \mathbf{U}_2\mathbf{V}_2^T + \mathbf{E}_2, \quad (\text{Equation 2})$$

where \mathbf{c} contains the direct effect of exposure on outcome, \mathbf{b} contains the effect sizes of DNAm levels on outcome, \mathbf{U}_2 are latent factors from a latent factor regression model, \mathbf{V}_2 are the corresponding loadings, and \mathbf{E}_2 is a matrix of errors. For each marker j , a significance value, P_y , is computed for the test of a null effect size for DNAm on outcome ($H_0: b_j = 0$).

Step 2. Identifying potential CpG mediators and aggregated mediator regions. The second step of HDMAX2 consists in combining the significance values P_x and P_y computed at each DNAm marker by using a new procedure called the max-squared test. The P -value for the max-squared (\max^2) test was computed as $P = \max(P_x, P_y)^2$. Like the Sobel test, the \max^2 test rejects the null-hypothesis that either the effect of exposure on DNAm or the effect of DNAm on outcome is null. The square in the formula warrants that the distribution of P -values is uniform when P_x and P_y are independent and uniformly distributed. In HDMAX2, the \max^2 test was first used in order to identify potential CpG mediators. A combination of P -values along the methylome was then performed to identify potential AMRs using *comb-p*, a method relying on the Stouffer-Liptak-Keckris correction that combines adjacent CpG P -values in sliding windows (Xu et al., 2016). We considered methylated regions including at least two markers at a maximum distance of 1,000 bp and significant at the 10% False Discovery Rate (FDR) level. The mean value of DNAm levels for CpGs located in AMRs was retained to summarize information on methylated regions.

Step 3. Quantifying indirect effects with single and multiple mediators. Mediation of exposure on the outcome was first assessed at the level of CpG markers, and then at the level of aggregated regions. For CpG and for AMRs, estimates of indirect effect sizes and the proportion of mediated effect were computed in the R package mediation (Imai et al.,

2010). For CpGs, the estimate of the indirect effect size for marker j was checked to be equivalent to the product of effect sizes, $a_j b_j$, computed in Eqs. (1) and (2). A novelty of HDMAX2 is to evaluate a cumulated indirect effect for all CpGs or AMRs identified in Step 3. The *overall indirect effect* (OIE) was estimated in a model including m mediator variables as follows

$$y_i = c x_i + \sum_{j=1}^m b_j m_{ij} + \sum_{k=1}^K v_k u_{ik}^{(2)} + \varepsilon_i$$

(Equation 3)

where m_{ij} represent methylation levels observed at CpG mediators or AMRs (or both of them), and the terms $u_{ik}^{(2)}$ correspond to the latent factor coordinates estimated in Step 1. The overall indirect effect was then computed as

$$OIE = \sum_{j=1}^m a_j b_j$$

(Equation 4)

where a_j represents the effect of exposure on methylation (Step 1). To account for correlation among mediators, the standard deviation of the OIE estimate was computed using a bootstrap approach (10,000 replicates).

2.2 Simulation studies

We performed simulations to compare the methods implemented in HDMAX2 with state-of-the-art approaches for EWAS in Step 1 and for mediation tests in Step 2.

Step-1 EWAS methods evaluated in simulations. In HDMAX2 Step 1, several latent factor estimation algorithms could be implemented for performing the EWAS. A preliminary study was performed to decide which of LFMM2, SVA and CATE was the best (using precision (1 - false discovery rate) and F1-score (harmonic mean of precision and power)) for our particular data set (Text S1, Figure S1). Then we performed generative simulations to compare methods using latent factors with those based on estimates of cell type composition. HDMAX2 was compared to two regression models including covariates obtained from RefFreeEWAS (Houseman et al., 2016) and refactor (Rahmani et al., 2016).

Step-2 mediation methods evaluated in simulations. We compared the \max^2 mediation test in Step 2 of HDMAX2 to methods based on direct application of Sobel tests and of univariate mediation analysis (Cardenas et al., 2019; Morales et al., 2016). Then we compared the \max^2 test to recent methods for high dimensional mediation: a multiple-testing procedure for high-dimensional mediation hypotheses similar to the max-squared test HDMT (Dai et al., 2020), a two-step familywise error rate procedure called ScreenMin (Djordjilović et al., 2020), an approach using familywise error rate and false discovery rate control when testing multiple mediators SBMH (Sampson et al., 2018), a linear regression model combined with an ANOVA (Tobi et al., 2018) and an approach using variable selection to reduce the number of mediators HIMA (Zhang et al., 2016). The last two approaches combined the first steps of HDMAX2 in a single step.

Mediation model simulations. The simulations were performed according to a generative model that reproduces the mediation pathways described in Eq. (1) and Eq. (2). Exposure and outcome (X and Y) and three confounding factors (\mathbf{U}) were simulated according to a multivariate Gaussian model. The percentage of variance of exposure and outcome explained by the confounding factors, and the correlation between those variables, were set at 10%. The variances of confounding factors were equal to one. The number of DNAm markers was equal to $m = 38,000$, approximately equal to the number of CpGs for a single chromosome in our empirical data, and the number of individuals was equal to $n = 500$. The vectors of effect sizes (\mathbf{a} for exposure and \mathbf{b} for outcome) were generated by setting a proportion of effect sizes to zero. Non-null effect sizes were sampled according to a standard Gaussian distribution. A residual error matrix \mathbf{E} was simulated by using a multivariate Gaussian distribution with means equal to zero and standard deviations of one. In addition to the three confounding factors, six additional factors representing distinct cell types were added in the simulation model. The proportion of cells from six different types were simulated by using a Dirichlet distribution. To consider values that are realistic with respect to our data analysis, the parameters of the Dirichlet distribution were equal to the proportions of each cell type estimated on the EDEN placental DNAm data (described afterwards). A matrix of DNAm markers was built using Eq 1 and Eq 2 with 3 parameters: the mean of non-null effect size for exposure (X) on methylation \mathbf{M} ($a = 0.2, 0.4$), the mean of non-null effect size for \mathbf{M} on outcome ($b = 0.2, 0.4$), and the number of causal markers (equal to 8, 16 or 32). For each set of parameters, 200 simulations were carried out. For each method tested, a subset of hits with a level of $FDR = 5\%$ was selected as potential mediators (Benjamini and Hochberg, 1995). For each list of hits, we computed precision ($1 - FDR$), sensitivity (power), and the harmonic mean of precision and sensitivity (F1-score). The highest value of

an F1-score is one, if precision and sensitivity are maximal, and the lowest value is zero, if either the precision or the sensitivity is null.

2.3 Mediation analyses of prenatal exposure to maternal smoking on birth weight and gestational age

Study population. Our analysis included participants of the EDEN Mother-Child Cohort enrolled in the university hospitals of Nancy and Poitiers, France, between 2003 and 2006 (Abraham et al., 2018; Heude et al., 2016). Lifestyle, demographic, and clinical data were collected by questionnaires and interviews during pregnancy and after delivery.

DNAm measurements. DNAm was measured from DNA extracted from 668 placental samples. The Illumina's Infinium HumanMethylation450 BeadChip was used to assess the levels of methylation in samples following the manufacturer's instructions (Illumina, San Diego, CA, USA). Protocols for placental DNA extraction and DNAm processing are detailed in (Jedynak et al., 2021). Briefly, DNAm was normalized using the beta-mixture quantile (BMIQ) method to ultimately obtain beta-methylation levels for 379,904 CpG probed CpG sites (Teschendorff et al., 2013).

Maternal smoking (MS), birth weight (BW) and gestational age (GA). We excluded preterm deliveries (n=28), women who reported quitting smoking before pregnancy (n=70) and women whose smoking status was unknown (n=100), leaving 470 women included in our analyses. Birth weight was extracted from medical records. Prenatal maternal cigarette smoking was collected by questionnaires during prenatal and postpartum clinical examinations. *Non-smokers* were defined as women who did not smoke during the 3 months before and during pregnancy (359 non-smokers). *Smokers* were defined as women smoking more than one cigarette per day throughout the duration of the pregnancy (111 smokers). All smokers during pregnancy also smoked during the 3 months before pregnancy. *Gestational age* was defined as gestational age at birth.

Mediation analyses. We hypothesized that maternal smoking during pregnancy could induce modifications of placental DNAm that result in changes in GA or in BW. To this aim, we investigated the causal relationships between MS, placental DNAm and each pregnancy outcome. MS was encoded as a categorical variable and the outcomes were encoded as continuous variables. In order to identify mediators of the exposure-outcome relationship, we

used the HDMAX2 approach to evaluate DNAm CpG mediators first, and then to identify AMRs.

In HDMAX2 regression models, adjustment factors included child sex, parity (0, 1, ≥ 2 children; categorical covariate), maternal age at end of education (≤ 18 , 19–20, 21–22, 23–24, ≥ 25 years; categorical covariate), season of conception (categorical covariate), study center, maternal body mass index (BMI) before pregnancy, maternal age at delivery, batch, plate, and chip technical factors related to DNAm measurements (categorical covariates). We relied on the principal component analysis of the DNAm matrix to include 6 latent factors in the HDMAX2 regression models (Figure S2). This number was consistent with the 6 factors selected in a previous work to represent the cell-types using the Refree algorithm (Rousseaux et al., 2019). FDR-corrected P -values were calculated for the 379,904 CpGs using the local FDR algorithm in *fdrtool* (Strimmer, 2008). Calibration of the \max^2 test P -values was evaluated through a direct examination of the histogram of P -values. The local FDR parameter (η_0) was computed to evaluate the proportion of null-hypothesis among the 379,904 tests. This proportion was estimated at $\eta_0 = 99.8$ – 99.9% , suggesting that an FDR level of 5% would be overly conservative (Figure S3). To agree with the value of η_0 , candidate CpGs were selected at FDR levels $< 10\%$, corresponding to adjusted $P < 9.03 \times 10^{-6}$ for BW and to adjusted $P < 3.27 \times 10^{-6}$ for GA. Results obtained after considering FDR levels $< 20\%$ are also reported.

Chained mediation of maternal smoking on birth weight. To better understand the causal pathways involving (six) genic regions that mediate the effect of MS both on GA and on BW, we considered GA as having reverse causality on DNAm levels. To assess reverse causality, we evaluated the indirect effects of targeted AMRs in a mediation analysis of GA on BW. For AMRs having a significant mediation P -value, each indirect effect and an overall indirect effect were computed from the above-described procedures.

Bioinformatic analyses. Promoter and enhancer regions were obtained from Illumina chip annotations. Gene annotations were obtained using the *FDb.InfiniumMethylation.hg19* package (Triche, 2014). Placental gene expression of annotated genes was compared to their gene expression in other tissues according to the Expression Atlas database (Papatheodorou et al., 2020). For every gene, Chauvenet's criterion was used to decide whether the gene was outlier for placental expression compared to other tissues. Functional annotation was made from the KEGG and the Gene Annotation databases (Kanehisa et al., 2021).

3. Results

3.1 Simulations

HDMAX2 was compared to several recent combinations of methods for multidimensional mediation analysis using simulation experiments. First, we compared the performances of latent factor models to other regression methods in estimating the association between exposure, DNAm levels, and outcome (Step 1 of HDMAX2). Then we compared the \max^2 mediation test to recently proposed tests (Step 2 of HDMAX2).

Performances of regression methods in step 1 of HDMAX2. A preliminary simulation study evaluated which of SVA, LFMM, or CATE provided the best estimation algorithm of latent factors for our empirical data set (Text S1). CATE and LFMM obtained better performance scores than SVA (Figure S1). LFMM runtimes were shorter than those of CATE, and LFMM performance scores were higher. Thus, we concluded that LFMM is the most appropriate for analysis of the EDEN cohort data, and we used it everywhere in subsequent assessments of HDMAX2. Using more general simulation experiments, we measured the relative performances of HDMAX2, that jointly estimates effect sizes and latent factors with LFMM, and methods based on a priori estimates of cell-type composition with RefFreeEWAS and ReFACTor (Figure 2). In all scenarios, the performances of the ReFACTor method were much lower than those of LFMM and RefFreeEWAS (Figure 2). For lower effect sizes of DNAm on outcome, LFMM and RefFreeEWAS reached close F1-scores, but LFMM obtained higher scores than RefFreeEWAS for higher effect sizes. All approaches obtained higher scores when more mediators were simulated, or when both the effect of exposure on DNAm and the effect of DNAm on outcome were higher. The results indicated that latent factor regression models outperformed methods that directly attempt to estimate cell-type composition from the DNAm data.

Figure 2. Relative performances of Step 1 approaches estimating latent factors versus inclusion of cell-type composition estimates. F1-score as a function of the number of mediators (16 or 32), effect size of exposure on DNAm ($X \rightarrow M$; low = 0.2, high = 0.4), and effect size of DNAm on outcome ($M \rightarrow Y$; low = 0.2, high = 0.4). Each simulation included 38,000 CpGs for 500 samples, with 6 cell types and 3 additional confounding factors.

Performances of mediation tests in the 2nd step of HDMAX2. Next, we compared HDMAX2 to five recent tests for high-dimensional mediation: HDMT, ScreenMin, SBMH, linear models combined with ANOVA (lm+anova), and HIMA (Figure 3). In every scenario, HDMAX2 and HDMT reached similar scores, and those approaches were the best ones overall. In the specific case of high DNAm on outcome effect sizes and low exposure on DNAm effect sizes, lm+anova obtained the best scores, immediately followed by HDMAX2 and HDMT. Lowest performances were obtained with ScreenMin, SBMH and HIMA. When both effect sizes were high, HIMA obtained the lowest performances. For low DNAm on outcome effect sizes, lm+anova and SBMH obtained the poorest performances. In addition, HDMAX2 outperformed mediation analyses combining EWAS with Sobel tests and with unidimensional mediation analyses repeated at each marker, especially when the number of mediators increased from 16 to 32 (Figure S4). Since the runtime was much shorter for HDMAX2 than for HDMT and for other approaches (Figure S5), HDMAX2 was used in our analyses on empirical data.

Figure 3: Relative performances of Step 2 multidimensional mediation methods. F1-score as a function of the number of mediators (16 or 32), effect size of exposure on DNAm (X->M; low = 0.2, high = 0.4), and effect size of DNAm on outcome (M->Y; low = 0.2, high = 0.4). Each simulation included 38,000 CpGs for 500 samples, with 6 cell types and 3 additional confounding factors.

3.2 Mediation of prenatal exposure to maternal smoking on birth weight and gestational age

Among 470 mother-infant pairs, mean maternal age at enrolment was 29 years (SD = 5 years), body mass index before pregnancy was 23 kg/m² (SD = 4.4 kg/m²) and 23.6% of women smoked during pregnancy (Table 1). Term birth weight (BW) ranged between 2010g and 4960g, with a mean of 3352 g +/- 435 g. Gestational age (GA) varied from 37 weeks to 42 weeks, with a mean of 40 weeks +/-1.20 week (8.4 days). Maternal smoking (MS) during pregnancy had a significant correlation with BW ($r = -0.16$, $P = 0.003$), but not with GA (Figure S6). BW and GA were significantly correlated in mother-infant pairs ($r = 0.31$, $P = 1.6 \times 10^{-12}$). After adjustment, the total effect of MS was 140 g lower BW (SD = 49.1 g, $P = 0.004$), and the total effect of MS was not significant for GA (effect size = 0.12, SD = 0.14, $P = 0.2434$).

Table 1. Characteristics of the study population (n = 470).

Characteristics	mean±SE	n (%)
Center		
Poitiers		189 (40)
Nancy		281 (60)
Sex of offspring		
Male		241 (51)
Female		229 (49)
Parity		
0		189 (40)
1		195 (41)
≥2		86 (18)
Maternal age at end of education (year)		
≤18		89 (20)
19-20		70 (15)
21-22		114 (24)
23-24		109 (23)
≥25		87 (18)
Season of conception		
January – March		100 (21)
April – June		103 (22)
July – September		130 (28)
October – December		137 (29)
Maternal smoking		
smoker		111 (24)
Non smoker		359 (76)
BMI (kg/m²)	23.0±4.4	
Maternal age (year)	29.4±5.0	
Birth weight	3352±435	
Gestational duration (weeks)	40±1.2	

BMI= pre-pregnancy Body Mass Index.

Mediation of maternal smoking on birth weight. A high dimensional mediation analysis of MS on BW was performed using placental DNAm data from the EDEN mother-child cohort. At an FDR level of 10% (5%), thirty-two (twenty) CpGs were identified as mediators of MS on BW (Table S1, Figure 4A, adjusted max-squared $P < 9.11 \times 10^{-6}$). Twenty CpGs were associated with a lower BW for the newborn (average ACME: -32.0 g, SD = 5.6 g; average proportion mediated (PM): 22.8%, SD = 4.0), and twelve CpGs were associated with a higher BW (average ACME: 32.6 g, SD = 10.3 g; average PM: 23.3 %, SD = 7.4) (Figure S7). The 32 CpGs were associated with an overall indirect effect corresponding to 40.3 g lower BW (SD = 51.3 g).

Examples of CpG mediators with the largest negative indirect effects include cg10624729 (adjusted $P = 5.15 \times 10^{-8}$), in *MIGA1* (Mitoguardin 1) a regulator of mitochondrial fusion, associated with 41 g lower BW, cg19406975 (adjusted $P = 9.27 \times 10^{-8}$), in *SH3BP5L* (SH3 Binding Domain Protein 5 Like) which functions as a guanine exchange factor, associated with 41 g lower BW, cg01686933 (adjusted $P = 6.98 \times 10^{-7}$), in *NECTIN1* (Nectin Cell Adhesion Molecule 1) which encodes an adhesion protein that plays a role in the

organization of epithelial and endothelial cells, associated with 41 g lower BW, and cg14502606 (adjusted $P = 1.04 \times 10^{-6}$), in *MLX* (MAX Dimerization Protein *MLX*) a transcription factor which plays a role in proliferation, determination and differentiation, associated with 38 g lower BW (Table S1).

Figure 4. High Dimensional Mediation analysis (HDMAX2) of maternal smoking on birth weight. A) Manhattan plot of $-\log(P\text{-values})$ obtained from HDMAX2 (CpGs). Gray bars without dots correspond to CpGs above the 20% FDR level. Gene names correspond to hits identified at the 10% FDR level (32 hits). Colored bars without gene names correspond to hits identified at the 20% FDR level (164 hits). B) Manhattan plot of $-\log_{10}(P\text{-values})$ for potential AMRs at the 10% FDR level (28 hits). C) Estimates of indirect effect (ACME) and proportions of mediated effect for confirmed AMRs (19 hits). For the A) and B) plots, the colors indicate the magnitude of indirect effects. For the C) plot, the colors indicate the distance to the gene (bp). Symbols on top of colored bars correspond to classification as enhancer, promoter or unknown. Overall indirect effect of AMRs: 52 g lower BW.

At an FDR level $< 20\%$, 164 mediators were discovered, including fifty-five CpGs within enhancer regions and twenty-six CpGs within promoter regions (Figure 4A). In comparison with the methylome, the list of mediators was enriched in hits corresponding to enhancer regions (33% of all hits, $P = 0.0003$, Fisher test, Figure S8A), and it was depleted in hits corresponding to promoter regions (15% of all hits, $P = 0.04$, Fisher test, Figure S8B). Several mediators were found in the body of a gene (109 hits), and some genes were hit more than once (*AJAP1*, *ESRP2*, *SH3BP2*, *SKI*, *SRSF5*, *VAV2* and *MLX*). We additionally performed mediation analyses for CpG cg27402634 (between *LINC00086* and *LEKR1*) and cg25585967 (*TRIO*) identified in (Morales et al. 2016), and for one CpG (cg11280108) in the HumanMethylation450 BeadChip which was among the seven CpGs identified in (Cardenas et al. 2019) from the EPIC chip. Although associations of DNAm with exposure to MS were significant for those CpGs (adjusted $P = 9.07 \times 10^{-14}$), none of those markers were mediators of MS on BW in our analysis ($q\text{-values} > 0.93$).

Regarding methylated regions, HDMAX2 detected twenty-eight potential AMRs, including four within enhancer regions, seven within promoter regions, and twenty within the body of a gene (FDR level $< 10\%$, Figure 4B). Nineteen AMRs were associated with statistically significant indirect effects ranging between 26.7 g lower BW and 33.0 g higher BW (Table S2). Twelve AMRs were associated with a lower BW (average ACME: -19.7 g, SD = 4.6; average PM: 14.0%, SD = 3.3%), and seven were associated with a higher BW (average ACME: 17.5 g, SD = 7.9; average PM: 12.5%, SD = 5.6%, Figure 4C). The 19 AMRs were associated with an overall indirect effect corresponding to 52 g lower BW (SD = 45 g). The overall indirect effect of both CpG mediators and AMRs was 44.5 g lower BW (SD = 60.7 g). The strongest evidence corresponded to AMR chr17:40,713,862-40,715,404 (adjusted $P =$

3.20×10^{-13}) in *COASY* (Coenzyme A Synthase) which plays an important role in numerous synthetic and degradative metabolic pathways in all organisms, associated with 26 g lower BW. This AMR was only 3kb close to another AMR, chr17:40,718,932-40,719,777 (adjusted $P = 9.37 \times 10^{-19}$), in *MLX*, which was associated with 27 g lower BW (Figure 4, Table S2 for a full list of AMRs).

Figure 5. Mediation analysis of maternal smoking on gestational age. A) Manhattan plot for CpG's $-\log(P\text{-values})$ obtained from HDMAX2. Gene names are displayed for hits obtained at the 10% FDR level (15 hits). Colored bars without gene names correspond to hits identified at the 20% FDR level (63 hits). Gray bars without dots correspond to CpGs above the 20% FDR level. B) $-\log_{10}(P\text{-values})$ for potential AMRs at 10% FDR level (31 hits). C) Estimates of indirect effects (ACME) and proportions of mediated effect for AMRs. Colors correspond to the intensity of the indirect effects. Symbols on top of colored bars correspond to categorization as enhancer, promoter, or unknown. Overall indirect effect of AMRs: 0.12 weeks lower GA.

Mediation of maternal smoking on gestational age. An independent mediation analysis was performed on the DNAm data in order to evaluate the indirect effects of MS on GA. At an FDR level $< 10\%$, fifteen CpGs (two CpGs at FDR level $< 5\%$) were identified as mediators of MS on GA (Table S3, Figure 5A, adjusted max-squared $P < 3.28 \times 10^{-6}$). The 15 CpGs were associated with a weak overall indirect effect corresponding to 0.28 week (2 days) lower GA (SD = 0.12) (Figure S9).

Examples of CpG mediators with the most negative effects include cg10298741 (adjusted $P = 4.82 \times 10^{-7}$), in *ZFH3* (Zinc Finger Homeobox 3) a transcription factor which regulates myogenic and neuronal differentiation, associated with 0.08 week lower GA, cg04908961 (adjusted $P = 9.19 \times 10^{-7}$), in *MIR17HG* (MiR-17-92a-1 Cluster Host Gene) a host gene for the MiR17-92 cluster, a group microRNAs (miRNAs) that may be involved in cell survival, proliferation, and differentiation, associated with 0.09 week lower GA, cg08402058 (adjusted $P = 1.04 \times 10^{-6}$), in *BLCAP* (Bladder cancer-associated protein) which reduces cell growth by stimulating apoptosis, associated with 0.09 week lower GA, (see Table S3 for a full list of CpG mediators). Ten CpGs were associated with a shorter GA (average indirect effect 0.09 week lower GA, SD=0.02; PM: 74%, SD = 14%), and five CpGs were associated with higher GA (average ACME: 0.09 week, SD = 0.01; PM: 71%, SD = 10%; Figure S9). At an FDR level $< 20\%$, sixty-three mediators were identified, including twenty-six hits within an enhancer region (Figure 5A). This subset of CpG mediators was enriched in hits corresponding to enhancer regions (33% of all hits, $P < 2.2 \times 10^{-16}$, Fisher test, Figure S8A).

The per-region analysis resulted in the detection of thirty-one potential AMRs, including eleven regions within enhancers, five within promoters, and twenty-six within the body of a

gene (Figure 5B, Table S4). Twenty-three AMRs were associated with small but statistically significant indirect effects ranging between -0.09 week and 0.10 week (none were associated with a significant mediated proportion). Five regions were associated with a lower GA (average ACME: -0.06 week ; SD = 0.01 ; average PM: 54.1%, SD = 13.1%), and eighteen regions were associated with a higher GA (average ACME : 0.06 week ; SD = 0.01; average PM: 49.4%, SD = 11.1%).

The 23 AMRs were associated with a weak overall indirect effect corresponding to 0.12 week (23 hours) shorter GA (SD = 0.11). The cumulative overall indirect effect of CpG mediators and AMRs was 0.09 week (16 hours) shorter GA (SD = 0.14). The largest negative indirect effects corresponded to AMR chr1:28,906,332 - 28,906,661 (adjusted $P = 7.89 \times 10^{-9}$) in *SNHG12*, (Small Nucleolar RNA Host Gene 12) an RNA gene that may promote tumorigenesis, associated with 0.09 week lower GA, chr20:36,148,579 - 36,149,354 (adjusted $P = 1.13 \times 10^{-10}$) in *BLCAP*, which encodes a protein that reduces cell growth by stimulating apoptosis, associated with 0.06 week lower GA, and chr17:40,714,100 - 40,714,374 (adjusted $P = 2.84 \times 10^{-6}$) in *COASY* associated with 0.06 week lower (Figure 5, Table S4 for a full list of AMRs).

Chained mediation of maternal smoking on birth weight through DNAm and GA. Six genes, *COASY*, *BLCAP*, *SKI*, *DECR1*, *ESRP2*, *PRRT1*, included AMRs that act as mediators both for MS on BW and for MS on GA. To better understand the causal pathways involving those genic regions, we tested the hypothesis that GA influences methylation levels in those regions, and estimated the indirect effects in a mediation analysis of GA on BW (Figure 6, Figure S10, Table S5). In this analysis, GA had significant indirect effects on BW for two of the six AMRs, in *COASY* (ACME = 6.9 g, mediation $P < 10^{-3}$), *BLCAP* (ACME = 5.1g, mediation $P = 0.01$). The two AMRs were associated with an overall indirect effect corresponding to 10 g higher BW (SD = 3.91).

In the genomic region surrounding the *COASY* gene (Figure S10), the AMRs were located in regions with low DNAm levels (Figure S11B), and MS decreased DNAm levels within AMRs (Figure S11C). The CpGs contained in AMRs mediated lower BW, and were among the most negative observed indirect effects (Figure S11D-E). In the genomic region surrounding the *BLCAP* gene (Figure S12), AMRs were located in highly methylated gene body areas (Figure S12B), and MS decreased DNAm levels within AMRs (Figure S12C). The CpGs contained in AMRs mediated lower BW, and again, they were among the most negative observed indirect effects (Figure S12D-E). Figure 6 provides a summary of the chained mediation analysis (Figure S13 for a summary of CpG mediation analysis).

Figure 6. Summary of mediation analysis of MS on GA and BW for AMRs. Nineteen AMRs mediate the relationship between MS and BW, with a total indirect effect corresponding to 52 g lower BW. Twenty-three AMRs mediate the relationship between MS and GA, with a total indirect effect corresponding to 0.12 week shorter GA. Two AMRs mediate the relationship between GA and BW with a total indirect effect corresponding to 10g higher BW. The color under each gene indicates the sign of the indirect effect. The color of segments indicates the direction of association (blue means negative, red means positive).

4. Discussion

Main contributions. High dimensional mediation holds promising results for deciphering molecular mechanisms underlying the association between exposure and outcomes. We present HDMAX2, a method combining estimates of latent factors in EWAS with max-squared tests for mediation, which also evaluates an overall mediated effect for CpG or AMR. We show that HDMAX2 outperforms state-of-art methods and recent approaches proposed to identify mediators in a high dimensional setting. HDMAX2 was applied to assess the indirect effects of exposure to MS on GA and BW in a study of 470 women from the EDEN mother-child cohort, and confirmed the important role played by placental DNAm in the causal pathway between maternal smoking during pregnancy and fetal growth outcomes (Nakamura et al., 2021). In addition to single CpG mediators, our analysis examined AMR and computed an overall indirect effect of all mediators considered simultaneously. The overall indirect effects of CpG and AMR were 44.5 g lower BW (32.1% of the total effect size) and 0.09 week lower GA (75% total effect size). These results support the hypothesis that the role of placental DNAm in the mediation of effect of exposure to MS on BW and on GA may be more polygenic than previously reported. In addition, a chained mediation analysis of MS on BW suggested the existence of reverse causal relationships for AMR located in the genes *COASY* and *BLCAP*, which mediate a proportion of the effect of MS on BW through an effect of GA on DNAm.

Simulation studies. The main improvements of HDMAX2 over existing mediation methods is the use of latent factor models for estimating hidden confounders in step 1, and the \max^2 test of mediation in step 2. The combination of latent factors and \max^2 tests proposed by the HDMAX2 approach were carefully evaluated with intensive simulations, and resulted in increased performances compared to five state-of-the-art methods evaluating multiple mediators (Dai et al., 2020; Djordjilović et al., 2020; Sampson et al., 2018; Tobi et al., 2018; Zhang et al., 2016). Latent factors increased statistical power compared to using a priori

estimates of cell-type proportions from reference-free methods (Houseman et al., 2016; Rahmani et al., 2016). The max-squared tests showed considerably better performances in comparison to the univariate mediation or Sobel test approaches, which were used in previous studies analyzing the role of placental DNAm data in the pathway between MS and BW (Cardenas et al., 2019; Morales et al., 2016). Using HDMAX2, none of the mediating CpGs identified using univariate mediation or Sobel test approaches (Morales et al. 2016 Cardenas et al. 2019) were mediators of MS on BW in our analysis (q -values > 0.93).

Mediation analysis of maternal smoking on birth weight. Previous studies have shown a possibly overestimated mediated effect of MS on BW, sometimes greater than the total effect size (Valeri et al., 2017). This is a limitation of univariate indirect effects estimated independently in a context of correlation between multiple mediators. In contrast, our approach estimated an overall indirect effect of the placental methylome representing 32% of the total effect size of MS on BW. Compared with previous placental DNAm mediation analyses of MS on BW (Cardenas et al., 2019; Morales et al., 2016), the magnitude of each mediator indirect effect size estimated in our cohort represented smaller part (less than 24% for AMRs) of the total effect size, and it was spread over more mediators suggesting that indirect effects are more polygenic than in previous estimates.

CpG mediators. HDMAX2 identified 32 CpG mediators of MS on BW, for which a majority (20/32) of effects represented a lower BW. The results provided evidence for an enrichment in enhancer regions and for a depletion in promoter regions among mediators, which agrees with conclusions from an association study between MS and placental DNAm in the EDEN cohort (Rousseaux et al., 2019). According to the Gene Ontology database, six mediators were located in genes linked to development or to the growth of tissues: cg24571086 in *FGFR2*, cg11362604 in *MEIS2*, cg00108098 in *SEMA5B*, cg10778780 in *CCK*, and cg20482145 in *MYH10* and cg07156115 in *AHR*. The genes *FGFR2* and *SEMA5B* are linked to the development of multicellular organisms and to the growth of developmental organs, *MEIS2* is linked to the development of the brain, eyes and pancreas, *CCK* is linked to neuron migration, and *AHR* is linked to the development of blood vessels.

Aggregated mediator regions. Evidence for increased polygenicity of placental DNAm mediation was confirmed by examination of AMRs, which are seen as more robust and more biologically meaningful than isolated differentially methylated CpGs (Svendson et al., 2016). HDMAX2 identified 19 AMRs of MS on BW, for which a majority of effects represented a lower BW (Figure 4C, Table S2). The most negative effects corresponded to AMRs in *COASY*, which plays an important role in numerous synthetic and degradative metabolic

pathways and in *MLX*, a transcription factor physically close to *COASY*, which is co-expressed in the placenta. Four regions were located in genes linked to tissue development or growth, in *FBN2* related to camera-type eye development, *ZFP42* to gonad development, *ESRP2* to fibroblast growth factor receptor signalling pathway, and *SKI* to roof of mouth development, olfactory bulb development, camera-type eye development, and skeletal muscle fiber development. The genes *FBN2* and *ZFP42* were over-expressed in the placenta compared to other tissues. Smoking-induced AMRs in *FBN2* and *ESRP2* were associated with higher BW, whereas AMRs in *ZFP42* and *SKI* were associated with lower BW. Looking more closely at the biology of mediators, we found a large number of them located in genes related to preeclampsia, a pregnancy complication of placental origin characterized by high blood pressure and protein in the urine, causing about a third of very premature births. Preeclampsia-related genes included *NECTIN1* (Ito et al., 2018), *AHR* (Wang et al., 2011), *FGFR2* (Marwa et al., 2016), *COASY* (Martin et al., 2015), *BLCAP* (Li et al., 2020), *SKI* (Martin et al., 2015), *AJAP1* (Yeung et al., 2016), *SH3BP5* (Kaartokallio et al., 2015). The over-representation of preeclampsia-related genes supports a pleiotropic effect of mediators, and highlights the difficulty of disentangling relationships between correlated outcomes.

Mediation analysis of maternal smoking on gestational age and reverse causality. Our results provided evidence that DNAm (CpG + AMR) mediates a relatively small total indirect effect of MS on GA, representing 0.09 week lower GA (15 hours). The largest negative effects corresponded to AMRs located in *SNHG12* and in *BLCAP* (Table S4). Six genes contained DMRs mediating both the effect of MS on BW and the effect of MS on GA. Two of those AMRs, located in *BLCAP* and *COASY*, had among the largest negative effects on both GA and BW. We found strong evidence that *BLCAP* and *COASY* were present in the causal pathway from GA to BW, but no evidence that they were present in the pathway from BW to GA. This result indicates that the corresponding AMRs in *BLCAP* and *COASY* may be involved in complex causal relationships, in which DNAm plays a role in the negative effect of MS on BW and GA amplified by a lower GA (Figure 6). The placenta exhibits a unique epigenetic profile, as it is one of the tissues with lower DNA methylation levels, which undergoes intense remodeling in early gestation, and dynamic changes with increase DNA methylation as gestation advances (Fuke et al., 2004; Novakovic et al., 2011). Whether placental DNAm influences GA or GA influences placental DNAm is uncertain. Our results suggest a bidirectional association between placental DNAm and GA, with a feedback loop from GA to BW through placental DNAm. However, a limitation of these results and of such investigation is the fact that GA and placental DNAm are co-occurring events.

Universally applicable framework for high dimensional mediating events. A large body of epigenetic research in perinatal health is dedicated to cord blood DNA methylation, although the placenta has attracted some attention (Cardenas et al., 2019; Everson et al., 2019; Morales et al., 2016). Indeed, the placenta supports both the health of the mother and the development of the foetus; it produces hormones, ensures immune-tolerance, and provide nutrients to the foetus and regulates the exchange of gases and wastes. The placenta is a key resource informing on the intra uterine environment and a highly relevant tissue to investigate within the DOHAD framework. Besides being associated with several prenatal exposures, placental DNA methylation is suggested to be a relevant proxy for neurodevelopmental outcomes (Jensen Peña et al., 2012; Kundakovic and Jaric, 2017; Lester and Marsit, 2018) and respiratory health (Chhabra et al., 2014) of the child and understanding the indirect effects of its DNAm modifications on such outcomes will be an important objective. Beyond the role of the placenta and DNA methylation, other tissues and omics markers are relevant to investigate in perinatal and more generally epidemiological studies. The HDMAX2 framework can be applied with other layer of mediators, basically any type of high-throughput data (i.e. gene expression data), or with data on any other tissue types.

5. Conclusions

We developed a novel algorithm for high-dimensional mediation, HDMAX2. Beyond our current application to placental DNAm data, HDMAX2 is applicable to a wide range of tissues and omics layers including genomics, transcriptomics, and other types of omics. HDMAX2 shows better performances on simulations and increased power compared to existing approaches. We showed the strength of HDMAX2 by applying it to characterize associations between exposure to maternal smoking during pregnancy and birth weight and gestational age at birth of the baby. The mediation analysis demonstrated a causal relationship between maternal smoking during pregnancy and those outcomes underpinning many more epigenetic regions than previously found, suggesting a polygenic architecture for the pathways. Not limited to single CpG markers, HDMAX2 is extended to identifying AMRs. AMRs provide more robust evidence than single epigenetic markers, and allowed the characterization of regions mediating effects of maternal smoking during pregnancy both on gestational age and birth weight, suggesting that placental DNAm is an important biological mechanism. We further showed the overall indirect effect accounting simultaneously for all

mediators identified as a plausible estimate of the mediated effect. AMRs located in *COASY* and *BLCAP* suggested reverse causality in the relationship between gestational and the methylome contributing to lower birth weight. Our study added several statistical improvements to high-dimensional mediation analyses, and revealed an unsuspected complexity of the causal relationships between maternal smoking during pregnancy and birth weight at the epigenome-wide level.

Abbreviations.

ACME	Average causal mediation effects
AMR	Aggregated Mediator Region
BMI	Body Mass Index
BMIQ	Beta Mixture Quantile
BW	Birth Weight of the baby
CATE	Confounder adjusted testing and estimation
CpG	Cytosine-phosphate-Guanine
DNAm	DNA methylation
DOHaD	Development Origins of Health and Diseases
EWAS	Epigenome Wide Association Studies
FDR	False Discovery Rate
GA	Gestational Age at delivery
HDMAX2	High-Dimensional Mediation Analysis with max-squared tests
HDMT	High-Dimensional Mediation Test
LFMM	Latent Factor Mixed Model
MS	Maternal Smoking during pregnancy
OIE	overall indirect effect
PM	Proportion Mediated
SVA	Surrogate Variable Analysis

6. Supplementary information

See supplemental materials including tables and figures

Software availability and requirements.

The method presented in this study is available in the R package HDMAX2 at <https://github.com/jumentib/HDMA> GNU and reusable under General Public License v3.0.

Acknowledgments

We thank Daniel Vaiman (INSERM U1016) for his help with lab experiments.

Author's contributions

BJ, CCB, ME performed the statistical analyses. JL, OF designed the study, wrote the manuscript and obtained funding. OF developed the statistical analysis plan. BH supervised the data collection and data management of the EDEN cohort, and provided guidance on the project. JT supervised the methylation and lab arrays. All authors read, revised and approved the final manuscript.

Funding

This work was supported by a grant from the French National Cancer Institute (INCa), the French Institute for Public Health Research (IreSP) (INCa_13641), and the French Agency for National Research (*ETAPE*, ANR-18-CE36-0005). BJ was partly supported by the Grenoble Alpes Data Institute, supported by the French National Research Agency under the Investissements d'Avenir program (ANR-15-IDEX-02), and by LabEx PERSYVAL Lab, ANR-11-LABX-0025-01. DNA methylation measurements were obtained thanks to grants from the Fondation de France (no. 2012-00031593 and 2012-00031617) and the French Agency for National Research (ANR-13-CESA-0011).

The EDEN mother-child study was supported by Foundation for medical research (FRM), National Agency for Research (ANR), National Institute for Research in Public health (IRESP), French Ministry of Health (DGS), French Ministry of Research, INSERM Bone and Joint Diseases National Research (PRO-A), and Human Nutrition National Research Programs, Nestlé, French National Institute for Population Health Surveillance (InVS), French National Institute for Health Education (INPES), the European Union FP7 programmes (FP7/2007-2013, HELIX, ESCAPE, ENRIECO, Medall projects), Diabetes National Research Program, French Agency for Environmental Health Safety (ANSES), Mutuelle Générale de l'Éducation Nationale (MGEN), French national agency for food security, Frenchspeaking association for the study of diabetes and metabolism (ALFEDIAM). We thank all the participants and members of the EDEN mother-child cohort study group.

Availability of data and materials

The EDEN datasets analyzed in the presented study are not publicly available as they are containing information that could compromise the research participant's privacy/consent. However, they are available from the corresponding author on reasonable request and with permission from the EDEN Steering Committee.

7. Declarations

Ethics approval and consent to participate

The EDEN cohort received approval from the ethics committee (CCPPRB) of Kremlin Bicêtre and from the French data privacy institution Commission Nationale de l'Informatique et des Libertés (CNIL). Written consent was obtained from the mother for herself and for the offspring.

Consent for publication

Not applicable.

Competing interests

The authors declare that they have no competing interests.

Author details

- (1) Université Grenoble-Alpes, Centre National de la Recherche Scientifique, Grenoble INP, TIMC CNRS UMR 5525, 38000 Grenoble, France.
- (2) Laboratory for Epigenetics and Environment, Centre National de Recherche en Genomique Humaine, CEA – Institut de Biologie François Jacob, University Paris Saclay, Evry, France
- (3) Université de Paris, CRESS, Inserm, INRAE, F-75004 Paris, France.
- (4) Inria Grenoble - Rhône-Alpes Inovalée 655 Avenue de l'Europe - CS 90051 38334 Montbonnot, France.
- (5) Université Grenoble-Alpes, INSERM, CNRS, Team of Environmental Epidemiology Applied to Development and Respiratory Health, Institute for Advanced Biosciences, 38000 Grenoble, France

References

Abraham, E., Rousseaux, S., Agier, L., Giorgis-Allemand, L., Tost, J., Galineau, J., Hulin, A., Siroux, V., Vaiman, D., Charles, M.-A., Heude, B., Forhan, A., Schwartz, J., Chuffart, F., Bourova-Flin,

- E., Khochbin, S., Slama, R., Lepeule, J., 2018. Pregnancy exposure to atmospheric pollution and meteorological conditions and placental DNA methylation. *Environ. Int.* 118, 334–347. <https://doi.org/10.1016/j.envint.2018.05.007>
- Agha, G., Hajj, H., Rifas-Shiman, S.L., Just, A.C., Hivert, M.-F., Burris, H.H., Lin, X., Litonjua, A.A., Oken, E., DeMeo, D.L., Gillman, M.W., Baccarelli, A.A., 2016. Birth weight-for-gestational age is associated with DNA methylation at birth and in childhood. *Clin. Epigenetics* 8, 118. <https://doi.org/10.1186/s13148-016-0285-3>
- Baron, R.M., Kenny, D.A., 1986. The Moderator-Mediator Variable Distinction in Social Psychological Research: Conceptual, Strategic, and Statistical Considerations 10.
- Benjamini, Y., Hochberg, Y., 1995. Controlling the False Discovery Rate: A Practical and Powerful Approach to Multiple Testing. *J. R. Stat. Soc. Ser. B Methodol.* 57, 289–300.
- Blum, M.G.B., Valeri, L., François, O., Cadiou, S., Siroux, V., Lepeule, J., Slama, R., 2020. Challenges Raised by Mediation Analysis in a High-Dimension Setting. *Environ. Health Perspect.* 128, 055001. <https://doi.org/10.1289/EHP6240>
- Cardenas, A., Lutz, S.M., Everson, T.M., Perron, P., Bouchard, L., Hivert, M.-F., 2019. Placental DNA Methylation Mediates the Association of Prenatal Maternal Smoking on Birth Weight. *Am. J. Epidemiol.* <https://doi.org/10.1093/aje/kwz184>
- Caye, K., Jumentier, B., Lepeule, J., François, O., 2019. LFMM 2: Fast and Accurate Inference of Gene-Environment Associations in Genome-Wide Studies. *Mol. Biol. Evol.* 36, 852–860. <https://doi.org/10.1093/molbev/msz008>
- Chhabra, D., Sharma, S., Kho, A.T., Gaedigk, R., Vyhlidal, C.A., Leeder, J.S., Morrow, J., Carey, V.J., Weiss, S.T., Tantisira, K.G., DeMeo, D.L., 2014. Fetal lung and placental methylation is associated with in utero nicotine exposure. *Epigenetics* 9, 1473–1484. <https://doi.org/10.4161/15592294.2014.971593>
- Dai, J.Y., Stanford, J.L., LeBlanc, M., 2020. A Multiple-Testing Procedure for High-Dimensional Mediation Hypotheses. *J. Am. Stat. Assoc.* 1–16. <https://doi.org/10.1080/01621459.2020.1765785>
- Djordjilović, V., Hemerik, J., Thoresen, M., 2020. On optimal two-stage testing of multiple mediators. *ArXiv200702844 Stat.*
- Djordjilović, V., Page, C.M., Gran, J.M., Nøst, T.H., Sandanger, T.M., Veierød, M.B., Thoresen, M., 2019. Global test for high-dimensional mediation: Testing groups of potential mediators. *Stat. Med. sim.* 8199. <https://doi.org/10.1002/sim.8199>
- Everson, T.M., Vives-Usano, M., Seyve, E., Cardenas, A., Lacasaña, M., Craig, J.M., Lesseur, C., Baker, E.R., Fernandez-Jimenez, N., Heude, B., Perron, P., González-Alzaga, B., Halliday, J., Deysenroth, M.A., Karagas, M.R., Íñiguez, C., Bouchard, L., Carmona-Sáez, P., Loke, Y.J., Hao, K., Belmonte, T., Charles, M.A., Martorell-Marugán, J., Muggli, E., Chen, J., Fernández, M.F., Tost, J., Gómez-Martín, A., London, S.J., Sunyer, J., Marsit, C.J., Lepeule, J., Hivert, M.-F., Bustamante, M., 2019. Placental DNA methylation signatures of maternal smoking during pregnancy and potential impacts on fetal growth (preprint). *Genomics.* <https://doi.org/10.1101/663567>
- Fuke, C., Shimabukuro, M., Petronis, A., Sugimoto, J., Oda, T., Miura, K., Miyazaki, T., Ogura, C., Okazaki, Y., Jinno, Y., 2004. Age related changes in 5-methylcytosine content in human peripheral leukocytes and placentas: an HPLC-based study. *Ann. Hum. Genet.* 68, 196–204. <https://doi.org/10.1046/j.1529-8817.2004.00081.x>
- Gao, Y., Yang, H., Fang, R., Zhang, Y., Goode, E.L., Cui, Y., 2019. Testing Mediation Effects in High-Dimensional Epigenetic Studies. *Front. Genet.* 10. <https://doi.org/10.3389/fgene.2019.01195>
- Heude, B., Forhan, A., Slama, R., Douhaud, L., Bedel, S., Saurel-Cubizolles, M.-J., Hankard, R., Thiebaugeorges, O., De Agostini, M., Annesi-Maesano, I., Kaminski, M., Charles, M.-A., 2016. Cohort Profile: The EDEN mother-child cohort on the prenatal and early postnatal determinants of

- child health and development. *Int. J. Epidemiol.* 45, 353–363. <https://doi.org/10.1093/ije/dyv151>
- Houseman, E.A., Kile, M.L., Christiani, D.C., Ince, T.A., Kelsey, K.T., Marsit, C.J., 2016. Reference-free deconvolution of DNA methylation data and mediation by cell composition effects. *BMC Bioinformatics* 17. <https://doi.org/10.1186/s12859-016-1140-4>
- Imai, K., Keele, L., Tingley, D., 2010. A general approach to causal mediation analysis. *Psychol. Methods* 15, 309–334. <https://doi.org/10.1037/a0020761>
- Ito, M., Nishizawa, H., Tsutsumi, M., Kato, A., Sakabe, Y., Noda, Y., Ohwaki, A., Miyazaki, J., Kato, T., Shiogama, K., Sekiya, T., Kurahashi, H., Fujii, T., 2018. Potential role for nectin-4 in the pathogenesis of pre-eclampsia: a molecular genetic study. *BMC Med. Genet.* 19, 166. <https://doi.org/10.1186/s12881-018-0681-y>
- Jedynak, P., Maitre, L., Guxens, M., Gützkow, K.B., Julvez, J., López-Vicente, M., Sunyer, J., Casas, M., Chatzi, L., Gražulevičienė, R., Kampouri, M., McEachan, R., Mon-Williams, M., Tamayo, I., Thomsen, C., Urquiza, J., Vafeiadi, M., Wright, J., Basagaña, X., Vrijheid, M., Philippat, C., 2021. Prenatal exposure to a wide range of environmental chemicals and child behaviour between 3 and 7 years of age – An exposome-based approach in 5 European cohorts. *Sci. Total Environ.* 763, 144115. <https://doi.org/10.1016/j.scitotenv.2020.144115>
- Jensen Peña, C., Monk, C., Champagne, F.A., 2012. Epigenetic effects of prenatal stress on 11 β -hydroxysteroid dehydrogenase-2 in the placenta and fetal brain. *PloS One* 7, e39791. <https://doi.org/10.1371/journal.pone.0039791>
- Joubert, B.R., Håberg, S.E., Nilsen, R.M., Wang, X., Vollset, S.E., Murphy, S.K., Huang, Z., Hoyo, C., Midttun, Ø., Cupul-Uicab, L.A., Ueland, P.M., Wu, M.C., Nystad, W., Bell, D.A., Peddada, S.D., London, S.J., 2012. 450K Epigenome-Wide Scan Identifies Differential DNA Methylation in Newborns Related to Maternal Smoking during Pregnancy. *Environ. Health Perspect.* 120, 1425–1431. <https://doi.org/10.1289/ehp.1205412>
- Jumentier, B., Caye, K., Heude, B., Lepeule, J., François, O., 2020. Sparse latent factor regression models for genome-wide and epigenome-wide association studies. <https://doi.org/10.1101/2020.02.07.938381>
- Kaartokallio, T., Cervera, A., Kyllönen, A., Laivuori, K., Kere, J., Laivuori, H., FINNPEC Core Investigator Group, 2015. Gene expression profiling of pre-eclamptic placentae by RNA sequencing. *Sci. Rep.* 5, 14107. <https://doi.org/10.1038/srep14107>
- Kanehisa, M., Furumichi, M., Sato, Y., Ishiguro-Watanabe, M., Tanabe, M., 2021. KEGG: integrating viruses and cellular organisms. *Nucleic Acids Res.* 49, D545–D551. <https://doi.org/10.1093/nar/gkaa970>
- Kundakovic, M., Jaric, I., 2017. The Epigenetic Link between Prenatal Adverse Environments and Neurodevelopmental Disorders. *Genes* 8, E104. <https://doi.org/10.3390/genes8030104>
- Küpers, L.K., Xu, X., Jankipersadsing, S.A., Vaez, A., la Bastide-van Gemert, S., Scholtens, S., Nolte, I.M., Richmond, R.C., Relton, C.L., Felix, J.F., Duijts, L., van Meurs, J.B., Tiemeier, H., Jaddoe, V.W., Wang, X., Corpeleijn, E., Snieder, H., 2015. DNA methylation mediates the effect of maternal smoking during pregnancy on birthweight of the offspring. *Int. J. Epidemiol.* 44, 1224–1237. <https://doi.org/10.1093/ije/dyv048>
- Lee, S., Sun, W., Wright, F.A., Zou, F., 2017. An improved and explicit surrogate variable analysis procedure by coefficient adjustment. *Biometrika* 104, 303–316. <https://doi.org/10.1093/biomet/asx018>
- Leek, J.T., Storey, J.D., 2007. Capturing Heterogeneity in Gene Expression Studies by Surrogate Variable Analysis. *PLoS Genet.* 3, 12.
- Lester, B.M., Marsit, C.J., 2018. Epigenetic mechanisms in the placenta related to infant neurodevelopment. *Epigenomics* 10, 321–333. <https://doi.org/10.2217/epi-2016-0171>
- Li, Yingying, Cui, S., Shi, W., Yang, B., Yuan, Y., Yan, S., Li, Ying, Xu, Y., Zhang, Z., Linlin

- Zhang, null, 2020. Differential placental methylation in preeclampsia, preterm and term pregnancies. *Placenta* 93, 56–63. <https://doi.org/10.1016/j.placenta.2020.02.009>
- MacKinnon, D.P., Fairchild, A.J., Fritz, M.S., 2007. Mediation Analysis. *Annu. Rev. Psychol.* 58, 593–614. <https://doi.org/10.1146/annurev.psych.58.110405.085542>
- Martin, E., Ray, P.D., Smeester, L., Grace, M.R., Boggess, K., Fry, R.C., 2015. Epigenetics and Preeclampsia: Defining Functional Epimutations in the Preeclamptic Placenta Related to the TGF- β Pathway. *PLOS ONE* 10, e0141294. <https://doi.org/10.1371/journal.pone.0141294>
- Marwa, B.A.G., Raguema, N., Zitouni, H., Feten, H.B.A., Olfa, K., Elfeleh, R., Almawi, W., Mahjoub, T., 2016. FGF1 and FGF2 mutations in preeclampsia and related features. *Placenta* 43, 81–85. <https://doi.org/10.1016/j.placenta.2016.05.007>
- Morales, E., Vilahur, N., Salas, L.A., Motta, V., Fernandez, M.F., Murcia, M., Llop, S., Tardon, A., Fernandez-Tardon, G., Santa-Marina, L., Gallastegui, M., Bollati, V., Estivill, X., Olea, N., Sunyer, J., Bustamante, M., 2016. Genome-wide DNA methylation study in human placenta identifies novel loci associated with maternal smoking during pregnancy. *Int. J. Epidemiol.* 45, 1644–1655. <https://doi.org/10.1093/ije/dyw196>
- Nakamura, A., François, O., Lepeule, J., 2021. Epigenetic Alterations of Maternal Tobacco Smoking during Pregnancy: A Narrative Review. *Int. J. Environ. Res. Public Health* 18, 5083. <https://doi.org/10.3390/ijerph18105083>
- Novakovic, B., Yuen, R.K., Gordon, L., Penaherrera, M.S., Sharkey, A., Moffett, A., Craig, J.M., Robinson, W.P., Saffery, R., 2011. Evidence for widespread changes in promoter methylation profile in human placenta in response to increasing gestational age and environmental/stochastic factors. *BMC Genomics* 12, 529. <https://doi.org/10.1186/1471-2164-12-529>
- Papathodorou, I., Moreno, P., Manning, J., Fuentes, A.M.-P., George, N., Fexova, S., Fonseca, N.A., Füllgrabe, A., Green, M., Huang, N., Huerta, L., Iqbal, H., Jianu, M., Mohammed, S., Zhao, L., Jarnuczak, A.F., Jupp, S., Marioni, J., Meyer, K., Petryszak, R., Prada Medina, C.A., Talavera-López, C., Teichmann, S., Vizcaino, J.A., Brazma, A., 2020. Expression Atlas update: from tissues to single cells. *Nucleic Acids Res.* 48, D77–D83. <https://doi.org/10.1093/nar/gkz947>
- Rahmani, E., Zaitlen, N., Baran, Y., Eng, C., Hu, D., Galanter, J., Oh, S., Burchard, E.G., Eskin, E., Zou, J., Halperin, E., 2016. Sparse PCA corrects for cell type heterogeneity in epigenome-wide association studies. *Nat. Methods* 13, 443–445. <https://doi.org/10.1038/nmeth.3809>
- Rousseaux, S., Seyve, E., Chuffart, F., Bourova-Flin, E., Benmerad, M., Charles, M.-A., Forhan, A., Heude, B., Siroux, V., Slama, R., Tost, J., Vaiman, D., Khochbin, S., Lepeule, J., the EDEN mother-child cohort study group, 2019. Maternal exposure to cigarette smoking induces immediate and durable changes in placental DNA methylation affecting enhancer and imprinting control regions (preprint). *Genomics*. <https://doi.org/10.1101/852186>
- Sampson, J.N., Boca, S.M., Moore, S.C., Heller, R., 2018. FWER and FDR control when testing multiple mediators. *Bioinformatics* 34, 2418–2424. <https://doi.org/10.1093/bioinformatics/bty064>
- Sobel, M.E., 1982. Asymptotic Confidence Intervals for Indirect Effects in Structural Equation Models. *Sociol. Methodol.* 13, 290. <https://doi.org/10.2307/270723>
- Strimmer, K., 2008. A unified approach to false discovery rate estimation. *BMC Bioinformatics* 9, 303. <https://doi.org/10.1186/1471-2105-9-303>
- Svendsen, A.J., Gervin, K., Lyle, R., Christiansen, L., Kyvik, K., Junker, P., Nielsen, C., Houen, G., Tan, Q., 2016. Differentially Methylated DNA Regions in Monozygotic Twin Pairs Discordant for Rheumatoid Arthritis: An Epigenome-Wide Study. *Front. Immunol.* 7, 510. <https://doi.org/10.3389/fimmu.2016.00510>
- Teschendorff, A.E., Marabita, F., Lechner, M., Bartlett, T., Tegner, J., Gomez-Cabrero, D., Beck, S., 2013. A beta-mixture quantile normalization method for correcting probe design bias in Illumina Infinium 450 k DNA methylation data. *Bioinforma. Oxf. Engl.* 29, 189–196.

- <https://doi.org/10.1093/bioinformatics/bts680>
- Thornburg, K.L., Marshall, N., 2015. The placenta is the center of the chronic disease universe. *Am. J. Obstet. Gynecol.* 213, S14-20. <https://doi.org/10.1016/j.ajog.2015.08.030>
- Tobi, E.W., Slieker, R.C., Luijk, R., Dekkers, K.F., Stein, A.D., Xu, K.M., Biobank-based Integrative Omics Studies Consortium, Slagboom, P.E., van Zwet, E.W., Lumey, L.H., Heijmans, B.T., 2018. DNA methylation as a mediator of the association between prenatal adversity and risk factors for metabolic disease in adulthood. *Sci. Adv.* 4, eaao4364. <https://doi.org/10.1126/sciadv.aao4364>
- Triche, T.J., 2014. *FDb.InfiniumMethylation.hg19*: Annotation package for Illumina Infinium DNA methylation probes. R package version 2.2.0.
- Valeri, L., Reese, S.L., Zhao, S., Page, C.M., Nystad, W., Coull, B.A., London, S.J., 2017. Misclassified exposure in epigenetic mediation analyses. Does DNA methylation mediate effects of smoking on birthweight? *Epigenomics* 9, 253–265. <https://doi.org/10.2217/epi-2016-0145>
- VanderWeele, T., 2015. *Explanation in Causal Inference: Methods for Mediation and Interaction*. Oxford University Press.
- VanderWeele, T.J., 2016. Mediation Analysis: A Practitioner's Guide. *Annu. Rev. Public Health* 37, 17–32. <https://doi.org/10.1146/annurev-publhealth-032315-021402>
- Wang, J., Zhao, Q., Hastie, T., Owen, A.B., 2017. Confounder adjustment in multiple hypothesis testing. *Ann. Stat.* 45, 1863–1894. <https://doi.org/10.1214/16-AOS1511>
- Wang, K., Zhou, Q., He, Q., Tong, G., Zhao, Z., Duan, T., 2011. The possible role of AhR in the protective effects of cigarette smoke on preeclampsia. *Med. Hypotheses* 77, 872–874. <https://doi.org/10.1016/j.mehy.2011.07.061>
- Xu, R., Hong, X., Zhang, B., Huang, W., Hou, W., Wang, G., Wang, X., Igusa, T., Liang, L., Ji, H., 2021. DNA methylation mediates the effect of maternal smoking on offspring birthweight: a birth cohort study of multi-ethnic US mother–newborn pairs. *Clin. Epigenetics* 13, 47. <https://doi.org/10.1186/s13148-021-01032-6>
- Xu, Z., Niu, L., Li, L., Taylor, J.A., 2016. ENmix: a novel background correction method for Illumina HumanMethylation450 BeadChip. *Nucleic Acids Res.* 44, e20. <https://doi.org/10.1093/nar/gkv907>
- Yeung, K.R., Chiu, C.L., Pidsley, R., Makris, A., Hennessy, A., Lind, J.M., 2016. DNA methylation profiles in preeclampsia and healthy control placentas. *Am. J. Physiol. Heart Circ. Physiol.* 310, H1295-1303. <https://doi.org/10.1152/ajpheart.00958.2015>
- Zeng, P., Shao, Z., Zhou, X., 2021. Statistical methods for mediation analysis in the era of high-throughput genomics: Current successes and future challenges. *Comput. Struct. Biotechnol. J.* 19, 3209–3224. <https://doi.org/10.1016/j.csbj.2021.05.042>
- Zhang, H., Zheng, Y., Zhang, Z., Gao, T., Joyce, B., Yoon, G., Zhang, W., Schwartz, J., Just, A., Colicino, E., Vokonas, P., Zhao, L., Lv, J., Baccarelli, A., Hou, L., Liu, L., 2016. Estimating and testing high-dimensional mediation effects in epigenetic studies. *Bioinformatics* 32, 3150–3154. <https://doi.org/10.1093/bioinformatics/btw351>

Step 1

EWAS using latent factor models

- $M \sim X$
- $Y \sim M + X$

Step 2 - CpG

Mediation test with \max^2 and detection of **CpGs** by **FDR** control

Step 2 - AMR

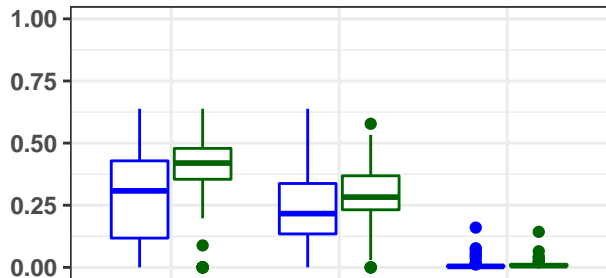
Detection of **AMRs** using the **comb-p** algorithm on the P-values of the \max^2 test

Step 3 - Estimate indirect effect

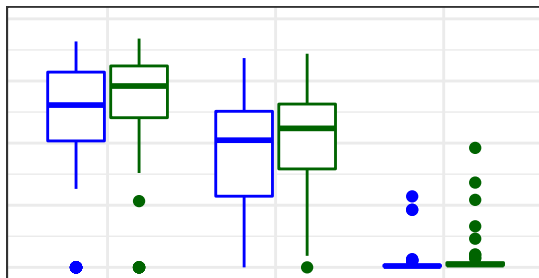
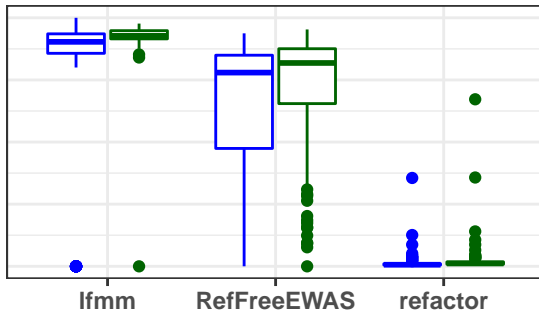
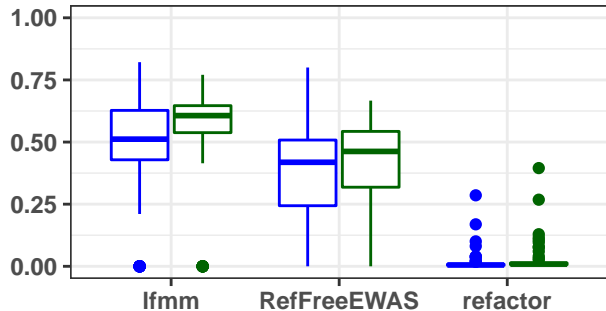
Estimation of the individual indirect effects of CpGs and AMRs via the **mediation** package

Estimation of the **overall indirect effects (OIE)** of CpGs, AMRs and both

Low Effect Size (M → Y)



High Effect Size (M → Y)

Low
Effect
Size
(X → M)Number of
mediatorsHigh
Effect
Size
(X → M)

ifmm

RefFreeEWAS

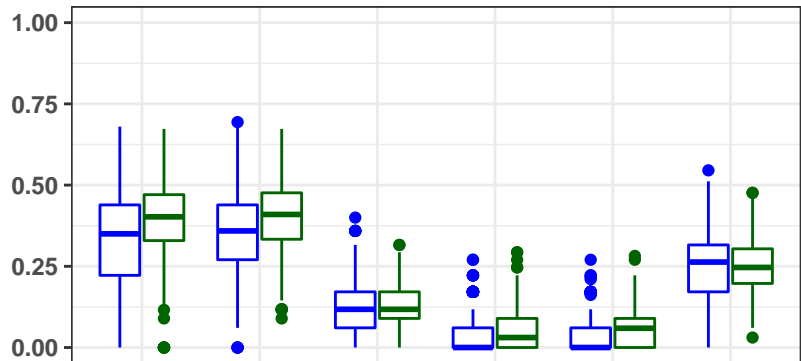
refactor

ifmm

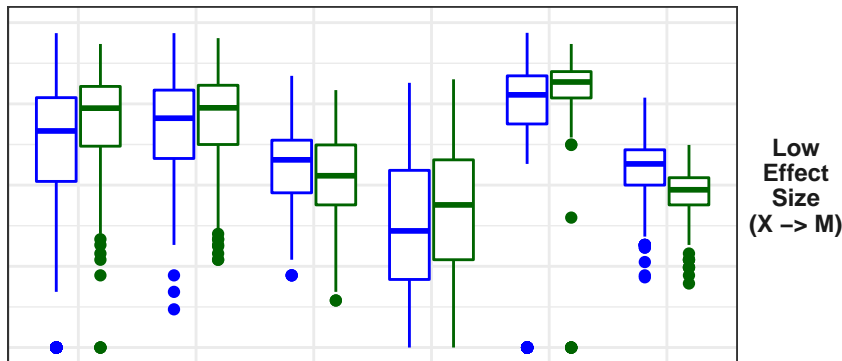
RefFreeEWAS

refactor

Low Effect Size (M → Y)

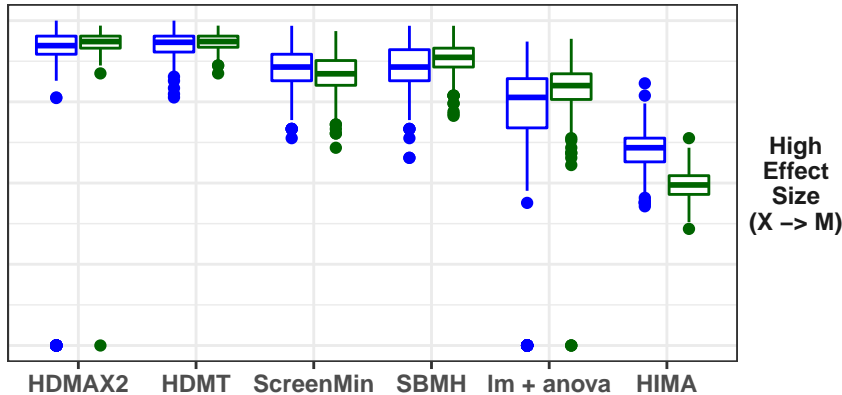
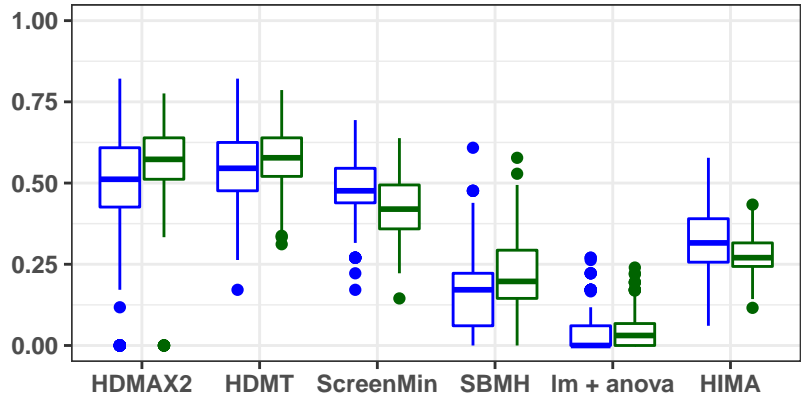


High Effect Size (M → Y)

Low
Effect
Size
(X → M)Number of
mediators

16

32

High
Effect
Size
(X → M)

HDMAX2

HDMT

ScreenMin

SBMH

Im + anova

HIMA

HDMAX2

HDMT

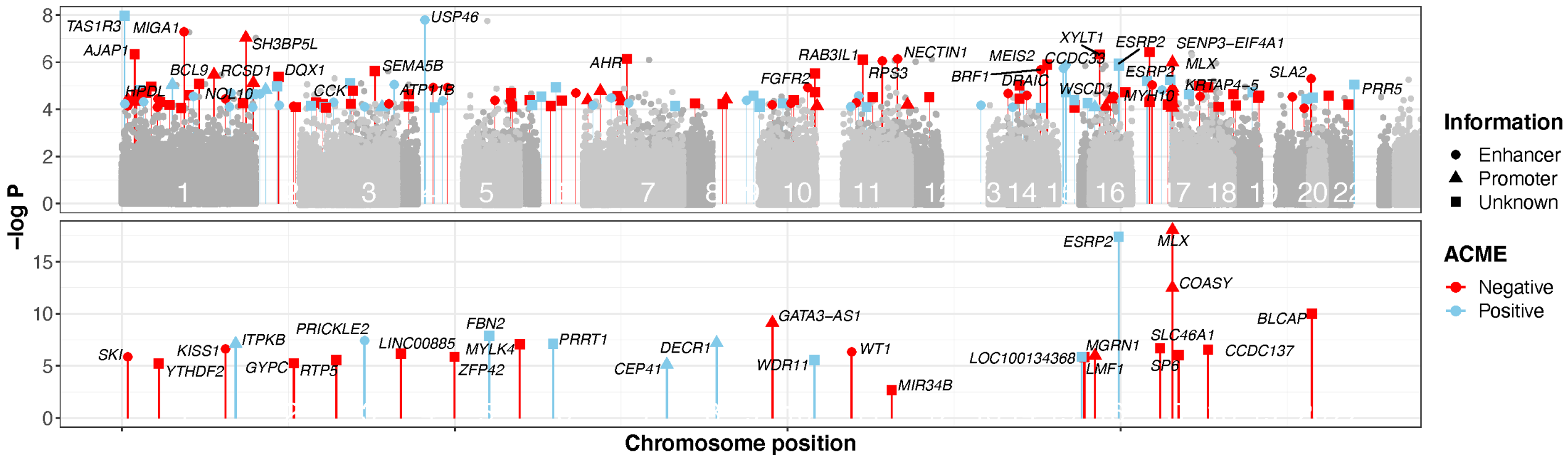
ScreenMin

SBMH

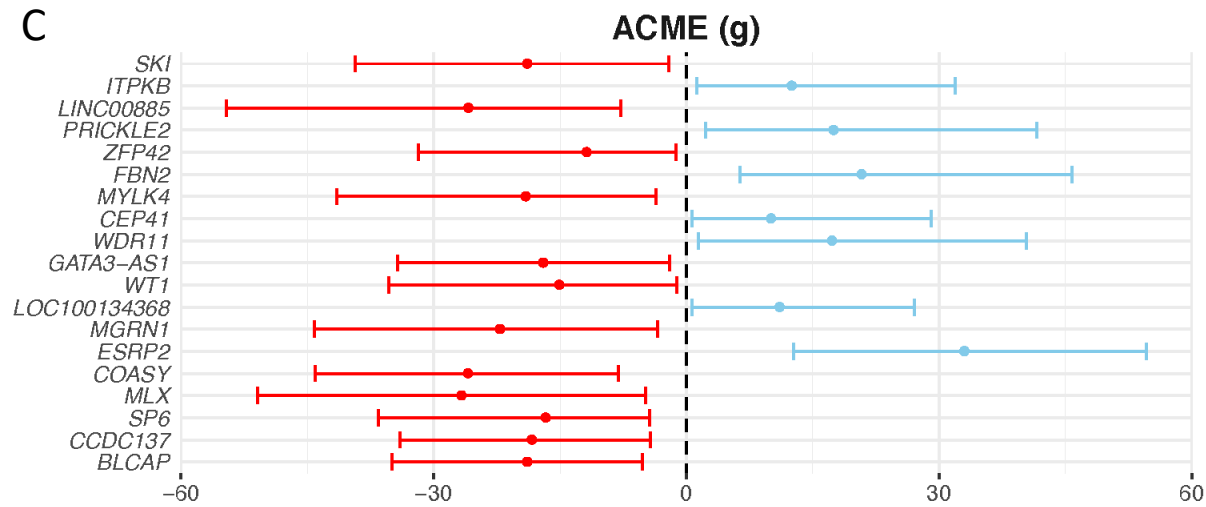
Im + anova

HIMA

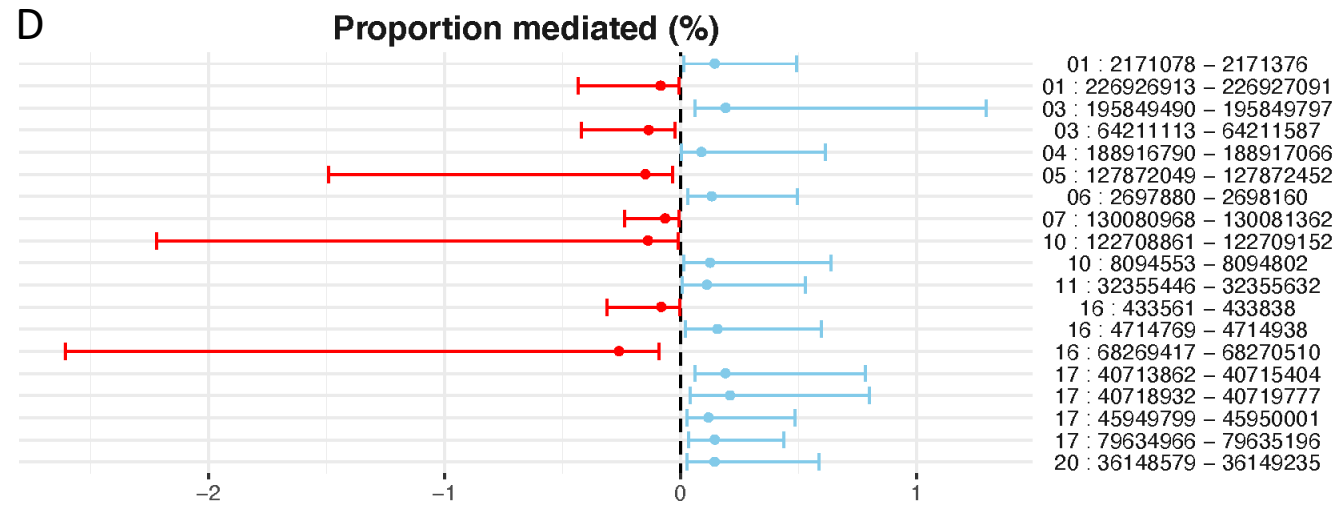
A - B



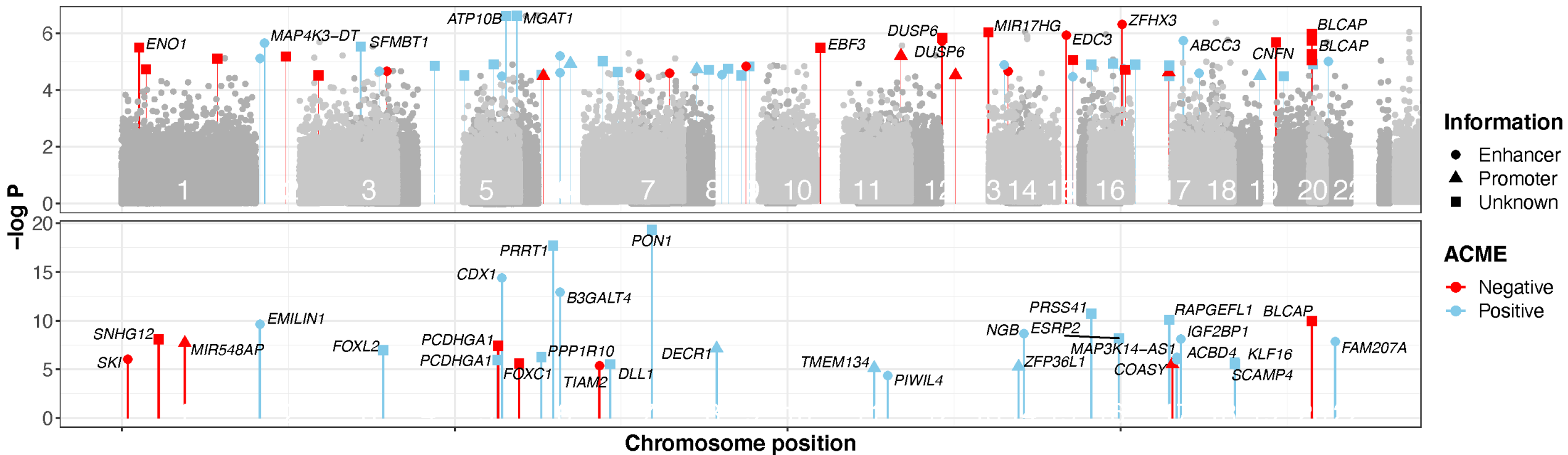
C



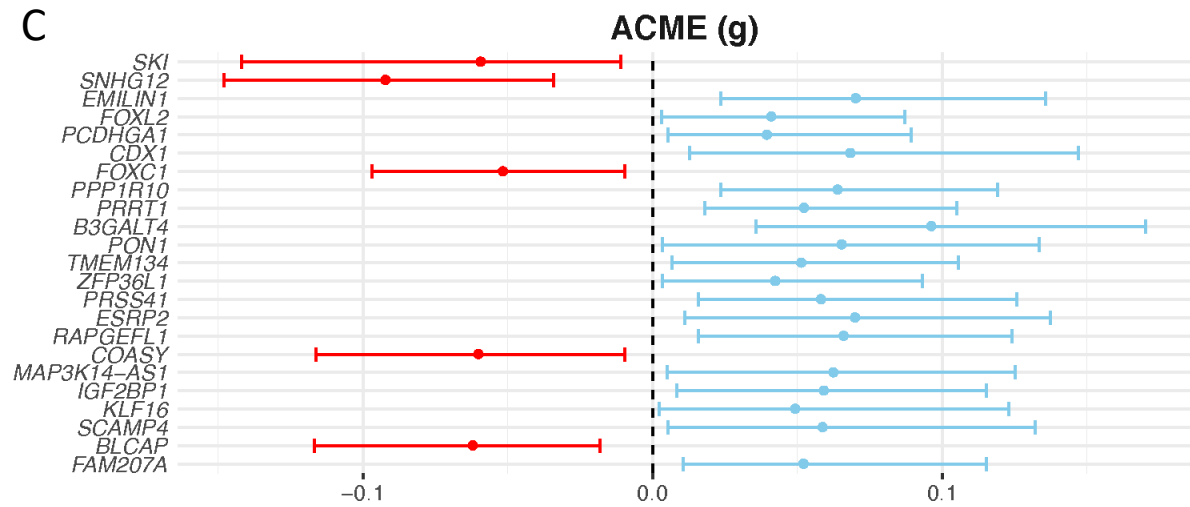
D



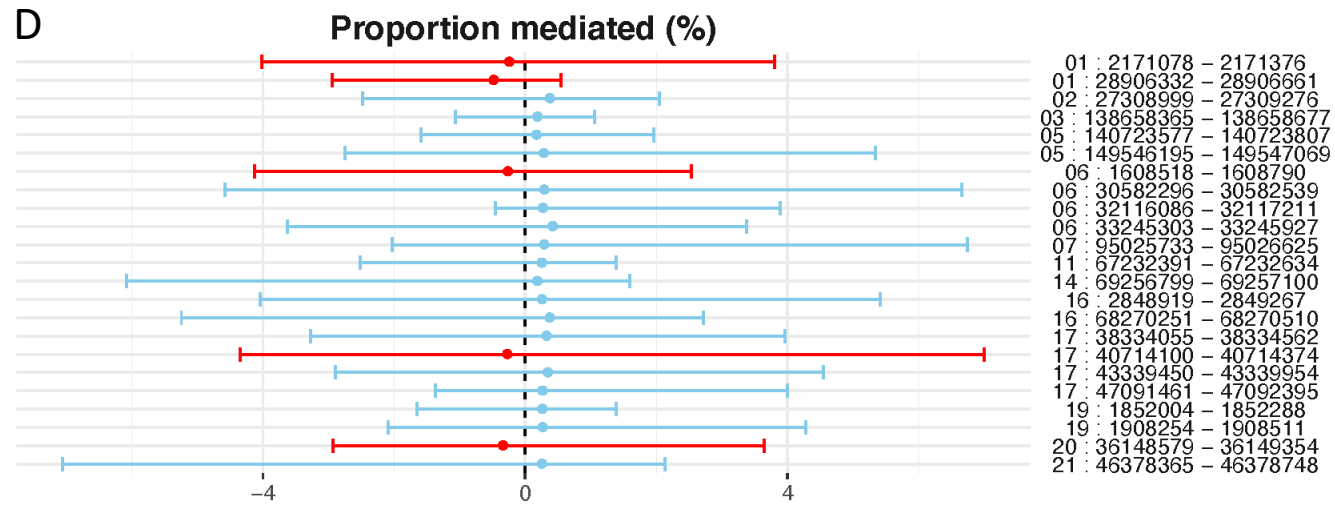
A - B



C



D



Relationships between MS and AMR and between AMR and BW or GA

— Negative
— Positive

medRxiv preprint doi: <https://doi.org/10.1101/2022.03.15.22272404>; this version posted March 16, 2022. The copyright holder for this preprint (which was not certified by peer review) is the author/funder, who has granted medRxiv a license to display the preprint in perpetuity. It is made available under a [CC-BY-NC-ND 4.0 International license](https://creativecommons.org/licenses/by-nc-nd/4.0/).

Indirect effect

■ Negative
■ Positive

