# QC EWAS

## QUALITY CONTROL FLOW FOR EPIGENOME-WIDE ASSOCIATION STUDIES
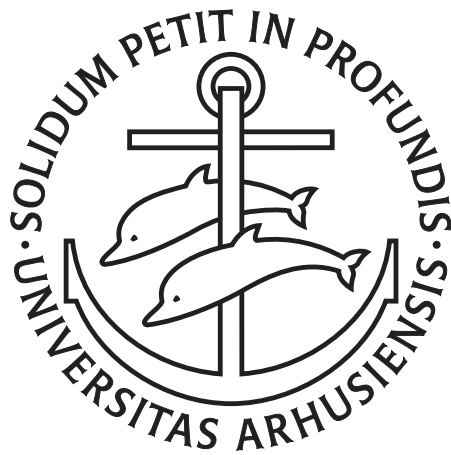
AINHOA SÁNCHEZ RAMÓN, AU726681

**AARHUS UNIVERSITET**

# QC EWAS

*Quality control flow for Epigenome-Wide Association Studies*

AINHOA

MSc in Bioinformatics

Bioinformatics Research Center
Faculty of Natural Sciences
Aarhus University

June 2024

# CONTENTS

## LIST OF FIGURES

## ACRONYMS AND ABBREVIATIONS

DNMTs   DNA methyltransferases

SAM     S-adenosylmethionine

5mC     5-methylcytosine

CGIs    CpG islands

TSSs    transcription start sites

iPSYCH  Integrative Psychiatric Research

ICC     intraclass correlation coefficient

DMSs    differentially methylated CpG sites

DPCRR   Danish Psychiatric Central Research Register

DNPR    Danish National Patient Register

ICD     World Health Organization International Classification of
        Disease

IDAT    Intensity Data

ASD     Autism spectrum disorder

MBR     Danish Medical Birth Register

ASD     autism spectrum disorder

nRBCs   nucleated red blood cells

QC      quality control

# 1

## ABSTRACT

Epigenetics is a essential field of study in understanding the complexities of gene regulation and expression, particularly in relation to environmental influences and their long-term effects on health and disease. The need for research investigating DNA methylation (DNAm) in clinical studies has greatly increased, leading to the evolution of new analytic methods to improve accuracy and reproducibility of the interpretation of results from these studies. Among the advancements in this field, the Illumina Infinium BeadChips have revolutionized large-scale epigenome-wide association studies (EWAS) in human populations. In most studies, the main objective of using DNAm analysis is to detect differences in methylation at CpG sites between phenotypic groups.

Despite its robustness this technology has limitations. Batch effects and confounding variables can introduce significant noise, potentially masking true biological signals. For this reason, rigorous quality control (QC) procedures are essential to ensure the reliability and validity of the data before any meaningful analysis can be performed. This thesis focuses on the application of comprehensive QC measures to the MINERvA dataset, which is inside the iPSYCH cohort. The iPSYCH cohort is a large and robust dataset that offers a unique opportunity to explore epigenetic mechanisms underlying complex traits in the Danish population.

We applied an extensive preprocessing into MINERvA to address various potential issues, including sample quality, probe reliability, and genetic ancestry. Initial quality assessments ensured that only samples meeting stringent criteria were included, thereby minimizing the risk of bias from low-quality data. Probe filtering was applied to retain only those probes that provided reliable measurements of DNAm. Ancestry filtering was employed to reduce confounding effects arising from population stratification, ensuring that observed associations were not artifacts of genetic background.

Although there are no standard methods for analyzing EWAS data, the application of established packages and methodologies already used in previous research made us sure of the dataset's integrity. Comprehensive data evaluation, including the use of statistical tests and visualizations, provided initial insights and confirmed that the data will be ready for subsequent analysis. The main goal of this is to amplify the statistical power of the next analyses of MINERvA,

reduce the risk of false positives, and support robust and reliable findings. Ensuring future reliable epigenetic studies of the MINERvA cohort. And increasing the potential for further research and clinical outcomes. By maintaining a high standard of data quality this research contributes to the field of epigenetics and its application in understanding and finding complex human diseases.