

Statistical Analysis Plan	PREDiCCT
Version No	1.0
Date Finalised	18/12/2020

4.1 Recruitment of patients, retention, and questionnaire completion

We will report number of patients confirmed eligible, included in the study, remaining in the study for two years follow-up (i.e. not withdrawing), and analysed. Where available, reasons will be given for participants formally withdrawing from the study.

The number and percentage of participants with baseline and monthly questionnaire data will be reported, at each time point and overall. Responses will be categorised as complete questionnaire; questionnaire submitted but <10% of questions unanswered; questionnaire submitted but 10-20% of questions unanswered; questionnaire submitted but 20% or more questions unanswered; or questionnaire missing. Number of patients with at most 10% of questions unanswered in at least 75% of the monthly follow-up questionnaires will be given.

4.2 Descriptive analysis

The total and average follow-up time will be computed across all patients, where follow-up time is measured as the length of time before the patient's final questionnaire is completed or patient withdrawal (whichever is later).

Frequency tables for the total number of hard and clinical (soft) flares reported will be constructed. A further descriptive summary will assess any association between the Covid-19 pandemic and flare event reporting. For the subset of participants whose follow-up spans the periods before and after commencement of lockdown in the UK, the hard and clinical (soft) flare tabulations will be further stratified into flares occurring up to and including 23 March 2020 and those occurring after 23 March 2020.

A Kaplan-Meier survival plot will be constructed showing time until first clinical flare and time to first hard flare on the same plot.

Descriptive analyses of demographic, clinical and social characteristics of study participants and their measured baseline exposures (see Table 1, section 4.4) will be calculated, split by clinical flare / no clinical flare and overall. These will be presented as mean (SD) for normally distributed variables, and median (interquartile range) for those not normally distributed. Binary, multinomial nominal and multinomial ordinal variables will be presented as number (percentage) in each category.

We will assess the differences between the clinical flare and no clinical flare populations using a t-test for normally distributed variables, the Mann–Whitney test for non-normally distributed continuous variables, and the Fisher's Exact test for differences in proportions.

4.3 Analysis of the Primary Outcome

The primary outcome is time to first clinical flare as assessed by the monthly questionnaires. A clinical flare is defined as the patient answering "no" to the question "Do you think your disease has been well controlled in the past 1 month?".

The primary exposure variables are:

1. Total animal protein intake (red meat, dairy, poultry, fish)
2. Dietary fibre (non-starch polysaccharides)
3. N-6 polyunsaturated fatty acids (PUFA)

Statistical Analysis Plan	PREDiCCT
Version No	1.0
Date Finalised	18/12/2020

4.5 Secondary outcome analyses

The same analysis as described in sections 4.3 and 4.4 will be repeated for the secondary outcome of hard clinical flare. This is defined as a clinical flare plus commencement of any new medication; altered dosing of existing medication for the treatment of IBD flare, with an increase in CRP (>5mg/L) and / or faecal calprotectin (>200mcg/g).

We will investigate the relationship of the exposures of interest with the secondary outcomes of (i) total number of clinical flares and (ii) total number of hard flares in the 24 months follow-up period using a negative binomial mixed effects regression analysis, where hospital site is the random effect and all others are fixed effects. Results will be expressed as rate ratios with 95% confidence intervals.

4.6 Validation and QC

A second statistician will separately program and check the primary analysis on the primary outcome (section 4.3), and all statistically significant effects on the primary outcome among secondary exposure variables. If the number of statistically significant effects is more than 20 then a random selection of 20 out of the total will be validated.

5. FFQ analysis plan

5.1 Descriptive analysis

Descriptive analysis of the dietary factors listed in Table 1 will be performed as described in section 4.2.

5.2 Primary analysis of the primary outcome

The primary dietary exposures of interest will be analysed with respect to the primary outcome as outlined in section 4.3. The exposures are:

1. Total animal protein intake (red meat, dairy, poultry, fish)
2. Dietary fibre (non-starch polysaccharides)
3. N-6 polyunsaturated fatty acids (PUFA)
4. Dietary emulsifiers (lecithin)

In addition to the covariates summarised in section 4.3, a categorical factor for season (as per the Met Office definition of spring [March/April/May], summer [June/July/August], autumn [September/October/November] and winter [December/January/February]) will be included as an adjustment variable in the analysis. Possible interactions between season and the associations between the primary dietary exposures and flare rate will also be investigated in exploratory analyses.

5.3 Secondary analysis of the primary outcome

The secondary exposures of interest N-3 polyunsaturated fatty acids (PUFA), sugar intake levels and starch intake levels will be analysed as in section 5.2.

Statistical Analysis Plan	PREDiCCt
Version No	1.0
Date Finalised	18/12/2020

5.4 Secondary outcome analyses

The secondary outcome analyses of section 4.5 will be applied to the FFQ exposures outlined in sections 5.2 and 5.3.

6. Microbiota data analysis plan

6.1 Quality control, mapping and quantification

Sequence data will be generated and undergo initial sequencing quality control. Quality control will include inspection of base quality scores, base composition and sequence duplication levels, and will identify poor quality sequencing lanes, runs or libraries to be removed. We will align the reads to the human reference genome, plus currently unassigned human contigs, using BWA-MEM to remove human contamination. We will then remove low quality reads using SolexaQA++ (Cox, 2010) on the default settings.

We will use MetaPhlAn2 (Truong, 2015) to quantify the relative abundance of different bacterial clades. We will use StrainPhlAn (Truong, 2017) and RAxML (Stamatakis, 2014) to construct phylogenetic trees of specific strains of interest across samples (see section 6.4). We will use HUMAnN2 (Abubucker, 2012) to calculate pathway abundances.

Counts for each microbiome measure (each species, genera, phyla and pathway) will be transformed using the centred log-ratio transformation, $\log(x/G(X))$, where $G(X)$ is the harmonic mean of read counts across all measures within that category (e.g. all species, for a species measure), as recommended by (Gloor, 2017). Zero counts will be imputed using the Zcompositions R package (Palarea-Albaladejo, 2015). Finally, we will remove microbiome measures that are inaccurately measured, due to a low abundance in a large number of samples (<2 reads observed in >50% of samples).

6.2 Descriptive analysis

We will generate a summary table of patient metagenomic characteristics, broken down by flaring and non-flaring patients. This will including total read count, % QC+, % human contamination, total number of observed species, total number of observed genes, and alpha diversity (measured by the inverse Simpson and Shannon indices).

We will test for associations between hospital site and microbiome measures, listing all phyla, genera, species or pathway that correlate significantly (Benjamini-Hochberg FDR < 0.05) with hospital site under a Kruskal-Wallis test.

We will generate two plots, each showing all samples on the first two principal coordinates of clades (species) and pathways (KO terms), coloured by hospital site. Principal coordinates will be calculated using the Bray–Curtis distance, using the R package vegan.

6.3 Primary analysis of the primary outcome