



HAL
open science

CroCoDeEL : accurate detection of cross-sample contamination in metagenomic data

Lindsay Goulet, Florian Plaza Oñate, Pauline Barbet, Edi Prifti, Eugeni Belda, Emmanuelle Le Chatelier, Guillaume Gautreau

► **To cite this version:**

Lindsay Goulet, Florian Plaza Oñate, Pauline Barbet, Edi Prifti, Eugeni Belda, et al.. CroCoDeEL : accurate detection of cross-sample contamination in metagenomic data: CroCoDeEL : accurate detection of cross-sample contamination in metagenomic data. 10th International Human Microbiome Consortium Congress, Jun 2024, Rome, Italy. . hal-04634008

HAL Id: hal-04634008

<https://hal.science/hal-04634008>

Submitted on 3 Jul 2024

HAL is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.



Lindsay Goulet¹, Florian Plaza Oñate¹, Pauline Barbet¹, Edi Prifti^{2,3}, Eugeni Belda^{2,3}, Emmanuelle Le Chatelier¹ and Guillaume Gautreau^{1,4}

Introduction

Background

The gut microbiota plays a crucial role in human health [1]. Metagenomic sequencing allows a deep characterization of microbial communities without prior organism isolation or culture.

Several massive sequencing projects are now on the launchpad as Le French Gut which aims to analyze 100 000 fecal samples to define the heterogeneity of healthy gut microbiota, the environmental and lifestyle factors impacting them, and their deviations seen in chronic diseases.

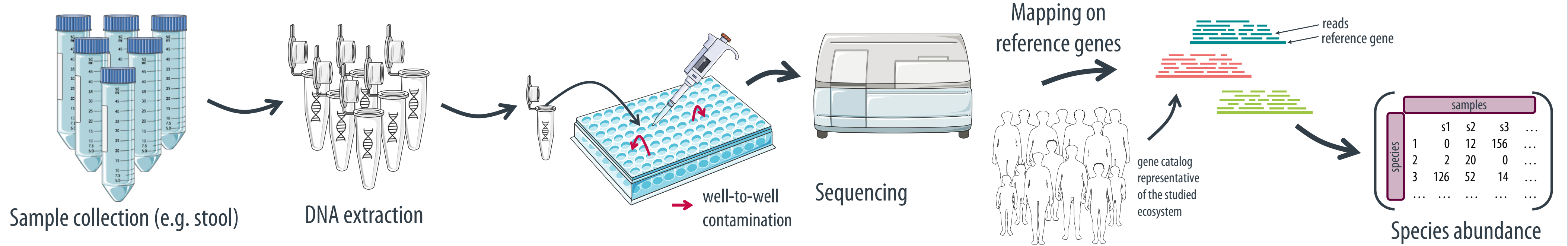
In this context, the detection of cross-sample contaminations is a crucial but time-consuming task. Checks must therefore be AI-assisted to ensure data quality at scale.

Cross-sample contamination

Contamination refers to the presence of DNA that does not originate from the biological sample under study. It can be due either to:

- DNA from an external source (environmental DNA [2] or lab reagents)
- DNA from another sample processed on the same plate (cross-sample/well-to-well contamination).

Cross-sample contamination occurs during wet lab steps (DNA extraction, sequencing library preparation).



Although cross-sample contamination is a common problem, it remains understudied. It can lead to biased results (i.e. overestimation of a diversity, false strain sharing events) and eventually to false conclusions if not detected and it is also a serious impediment to studies reproducibility.

We introduce CroCoDeEL, a tool based on a supervised pre-trained model to automatically detect cross-sample contamination. Contrary to state-of-the-art approaches [3][4], CroCoDeEL works with related samples that may naturally share strains (e.g.: mother/child), discriminates contamination sources from contaminated samples, estimates contamination rates and does not require costly negative controls.

Methods

How to detect cross-sample contamination ?

- Above a certain threshold, all the abundant species of the contamination source sample are present in the contaminated sample.
- A subset of these shared species have a proportional abundance between the two samples and form a contamination line (○). This line corresponds to species that were not present before the contamination event (○).
- The contamination line is used to detect contamination events, and the contamination rate can be estimated from relative abundance ratio of species that constitute it.

CroCoDeEL

Cross-sample Contamination Detection and Estimation of its Level

- Input -

species	s1	s2	s3	...
1	0	12	156	...
2	2	20	0	...
3	126	52	14	...

- Output -

source	target	rate	score
s1	s2	1%	0,8
s2	s3	2,5%	0,9

Step1: selecting species

Step2: linear modeling using RANSAC

Step3: extraction of 10 features

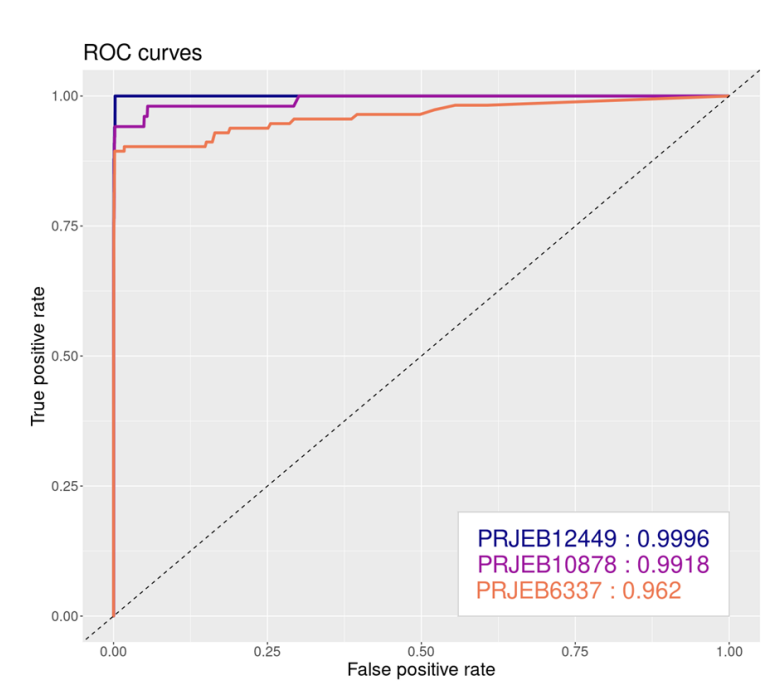
Step4: classification using reference random forest model (supervised learning)

Results

Performance

Performance was evaluated on 3 public metagenomics cohorts, curated by human experts :

	PRJEB12449	PRJEB10878	PRJEB6337
# of comparisons (# of plot×2)	11990	16256	55932
False Positive (FP)	39	38	101
False Negative (FN)	0	3	12
Accuracy	99.70%	99.70%	99.80%
Area Under the ROC Curve (AUC)	0.9996	0.9918	0.962
Runtime (seconds)	8 CPU: 133 16 CPU: 75	184 103	681 370
RAM consumption (MB)	8 CPU: 205 16 CPU: 206	222 206	213 214



Application

CroCoDeEL has been applied to fifteen public cohorts on colorectal cancer. Here are some non-exhaustive results :

project	paper	# samples	seq. depth (M reads) (mean ± sd)	# contaminated samples	contamination rate (Q1 med Q3)
PRJDB4176	YashidaS_2019	645	21.9 ± 6.3	5 (0.8%)	0.05% 0.08% 0.52%
PRJNA389927	HanniganGD_2018	84	2.8 ± 1.7	4 (4.8%)	5.57% 33.72% 64.94%
PRJEB27928	WirbelJ_2019	82	18.2 ± 8.1	39 (58.5%)	0.25% 0.35% 0.54%
PRJNA531273	GuptaA_2019	30	4.6 ± 2.1	17 (56.7%)	13.40% 36.33% 62.57%

Conclusion and future work

- We systematically found cross-sample contamination in the projects we curated but the proportion of contaminated samples and contamination rates varied significantly.
- These results show that this issue is widespread although some wet lab pipelines are more prone to cross-sample contamination.

- CroCoDeEL accurately detects cross-sample contamination in shotgun metagenomic data.
- CroCoDeEL is fast and works with limited computing resources.
- We believe that checking for sample cross-contamination should be a mandatory quality control step in metagenomics.
- CroCoDeEL can now be applied to your dataset, even without negative controls.
- Next efforts will focus on the automatic decontamination of samples, where possible.
- CroCoDeEL source code is available on GitHub : <https://github.com/metagenopolis/CroCoDeEL>
- CroCoDeEL can be easily installed with pip or conda.

