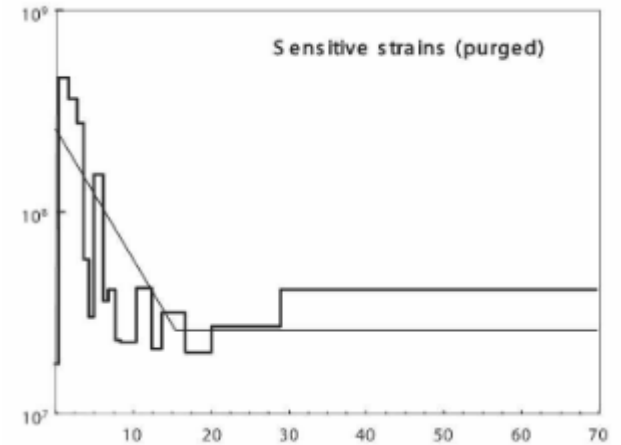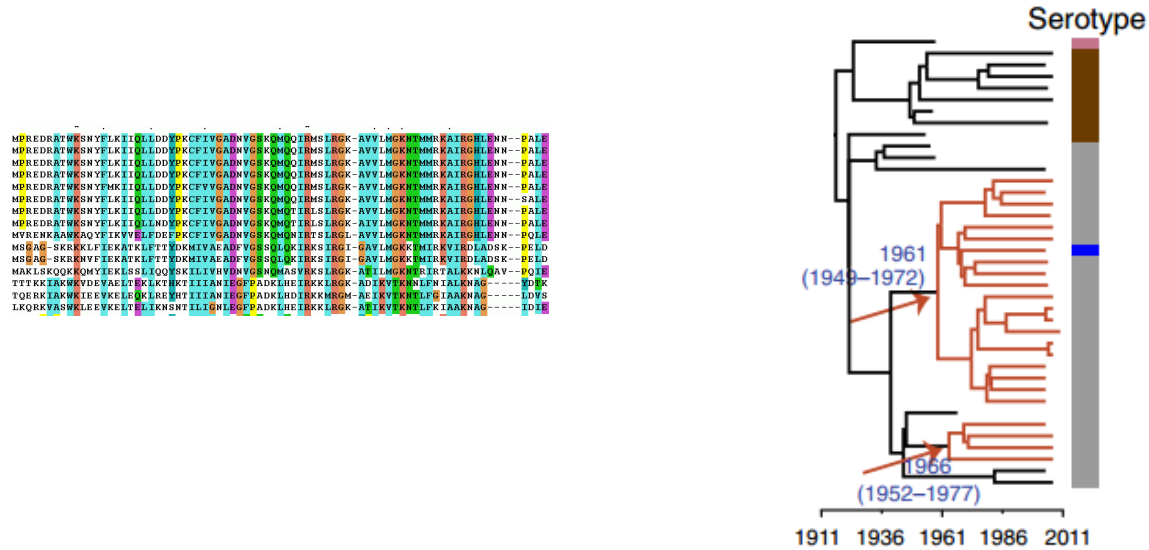# Inferring population sizes of bacterial populations
## a deep learning approach

### MLMicrobial Genomics - ECML -2022

Jean Cury, Théophile Sanchez, Erik Bray, Jazeps Medina-Tretmanis, Maria Avila-Arcos, Emilia Huerta-Sanchez, Guillaume Charpiat, and Flora Jay
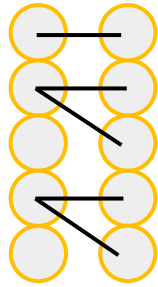
# Bacterial population genetics



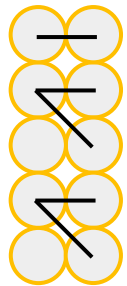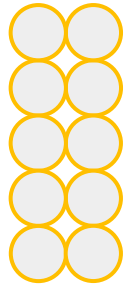⇒ Focus on population size inference

# Intuition

# Intuition

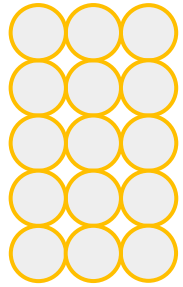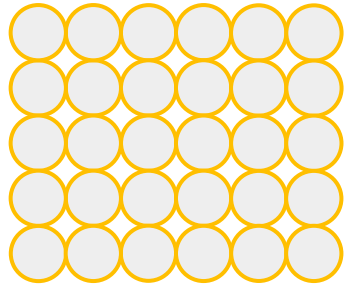Parental        Daughter cells

# Intuition

# Intuition

Generations →

# Intuition



Generations
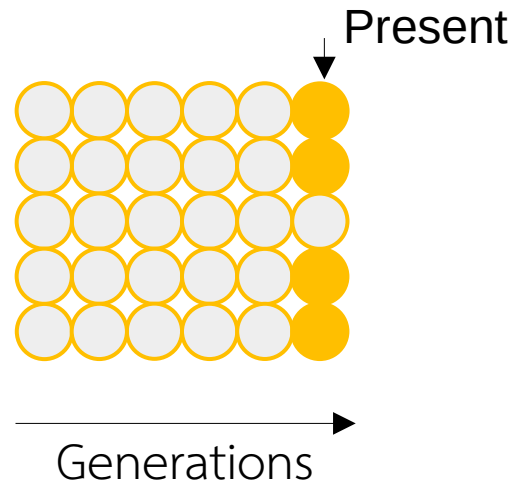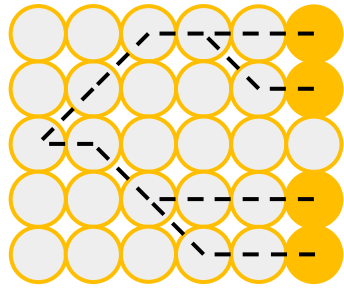
# Intuition



Generations →

# Intuition

Present



Generations

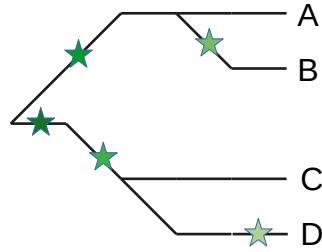# Intuition



Generations

# Intuition

# Intuition

# Intuition

# Intuition

# Intuition



Add selection sign

# Inference : Skyline Plot



- Using estimated coalescent time, it infers population size

- Non-parametric

- Does not require simulation

- Problem : does not work in bacteria

Lapierre et al. 2016

# Project

- End to end deep learning approach for bacterial popgen



- **Problem:** No ground truth data

- We need a population genetic simulator that is :

    – fast

    – Implement bacterial recombination (homologous HGT)

    – Demography, selection, etc..

# Input Data

MSA

```
ATGCGACAG
CTGCGTCGG
ATGAGTCAG
CTGCGTCAG
123456789
```

Observation

Population Size

Generations

# Input Data

MSA

SNP matrix

```
ATGCGACAG        000000000        0000
CTGCGTCGG        100001010        1011
ATGAGTCAG        000101000        0110
CTGCGTCAG        100001000        1010
123456789        123456789        1468
```

Observation

Population Size

Generations

# Input Data

MSA

SNP matrix

ATGCGACAG
**C**TGCG**TC**G**G**
ATG**A**G**T**CAG
**C**TGCG**T**CAG
123456789

000000000
**1**0000**101**0
000**1**0**1**000
**1**0000**1**000
123456789

0000
**1**0**11**
0**11**0
**1**0**1**0
1468

Observation

Population Size

Generations

Predictions

# Input Data

MSA

$x$         $y$

SNP matrix      ⟷    "Expansion"

```
ATGCGACAG          000000000          0000
CTGCGTCGG          100001010          1011
ATGAGTCAG          000101000          0110
CTGCGTCAG          100001000          1010
123456789          123456789          1468
```

Simulations

Population Size

Generations

Observation

Predictions

# What we simulate

### Decline

### Constant

### Expansion

Ne

Generations

5000x100 simulations

1000x100 simulations

5000x100 simulations

→ 50% with and 50% without <u>selection</u>

→ Variable parameters:

- **initial population size** (~Ne)
- mutation rate
- recombination rate (ratio r/m)
- coefficient of selection
- time of selection
- **time of demographical change**
- **strength of bottleneck/expansion**

→ Generated with a generalized Halton sequence

→ Fixed parameters:
- chromosome size
- mean size of gene conversion tracts
- Number of generations

Using SLiM, adapted for
Bacterial population
(Cury et al., 2021)

# Approach

- Use of **dnadna**, a package that help to reproduce, share and develop DL methods for population genetics

- Use of SPIDNA architecture
  - Invariant to permutation of individuals
  - Adaptive to input dimension
  - Good performance on human populations

- Add uncertainty estimation

Sanchez et al. 2022
Sanchez et al. 2020

# **dnadna** : Package for DL in population genetics



Package that allow:
- Development of network
- Reuse of someone else's network
- Reproduce training/prediction

→ Without coding skill (YAML)

https://gitlab.com/mlgenetics/dnadna

Sanchez et al. 2022

# Inference of demography

```
0000
1011
0110
1010
1468
```

400 first SNPs

SPIDNA

Population sizes at
21 time steps

# Inference of demography

Example with 100 simulations with the same set of parameters



Prediction on 1 simulation

Target value

# What about uncertainty ?

- DNN output a single value without notion of uncertainty :

    - Aleatoric : due to the underlying process that is intrinsically stochastic

    - Epistemic : Your sample is out of the distribution of the simulations



Out of distribution    Training distribution

# What about uncertainty ?

- DNN output a single value without notion of uncertainty :

  - Aleatoric : due to the underlying process that is intrinsically stochastic

  - Epistemic : Your sample is out of the distribution of the simulations

Use of Gaussian Negative Log Likelihood Loss to learn a gaussian with parameters μ and $\sigma^2$



Nix and Weigend, 1994
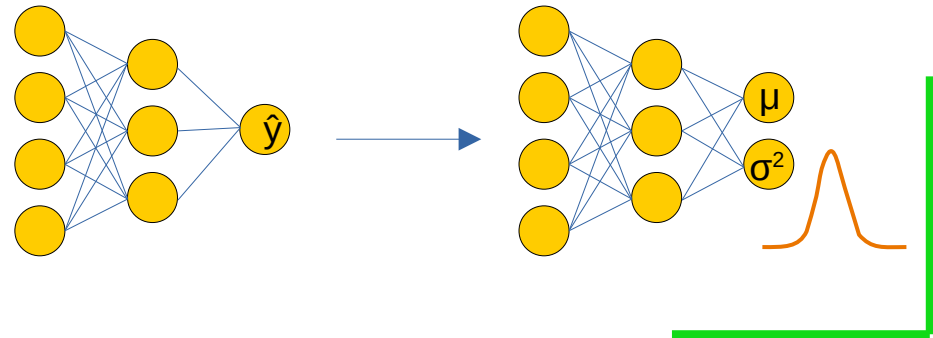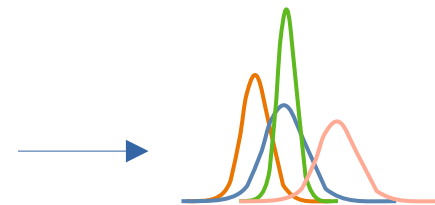Lakshminarayanan et al, 2017

# What about uncertainty ?

- DNN output a single value without notion of uncertainty :

  - ~~Aleatoric : due to the underlying process that is intrinsically stochastic~~

  - Epistemic : Your sample is out of the distribution of the simulations

Ensemble of Networks



Weighted mixture of Gaussian distribution

weights $\propto 1/\sigma^2$

Lakshminarayanan et al, 2017

# Uncertainty estimation



Population size

Without selection

With selection

50% High density Interval

Decline

90% High density Interval

Time

# Uncertainty estimation



Without selection

With selection

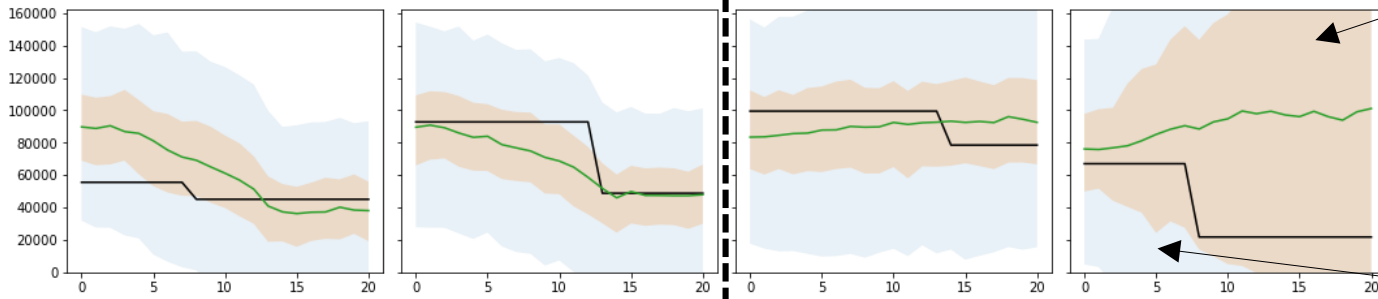50% High density Interval

90% High density Interval

Decline

Constant

Expansion

Population size

Time

# Error on Test set (the lower the better)



Bad predictions for Expansion

Good prediction otherwise

Except for ancient times where predictions follows the prior of the training set.

# Calibration of the Gaussian mixture

# ancient DNA

- Increasing amount of aDNA sequenced as technology improves
- Can help palaeontologist / historian understand  distant past
- Problem : low quality of sequences
    - Due to degradation of DNA
    - Higher rate of sequencing error
    - Poor coverage (small amount of DNA)

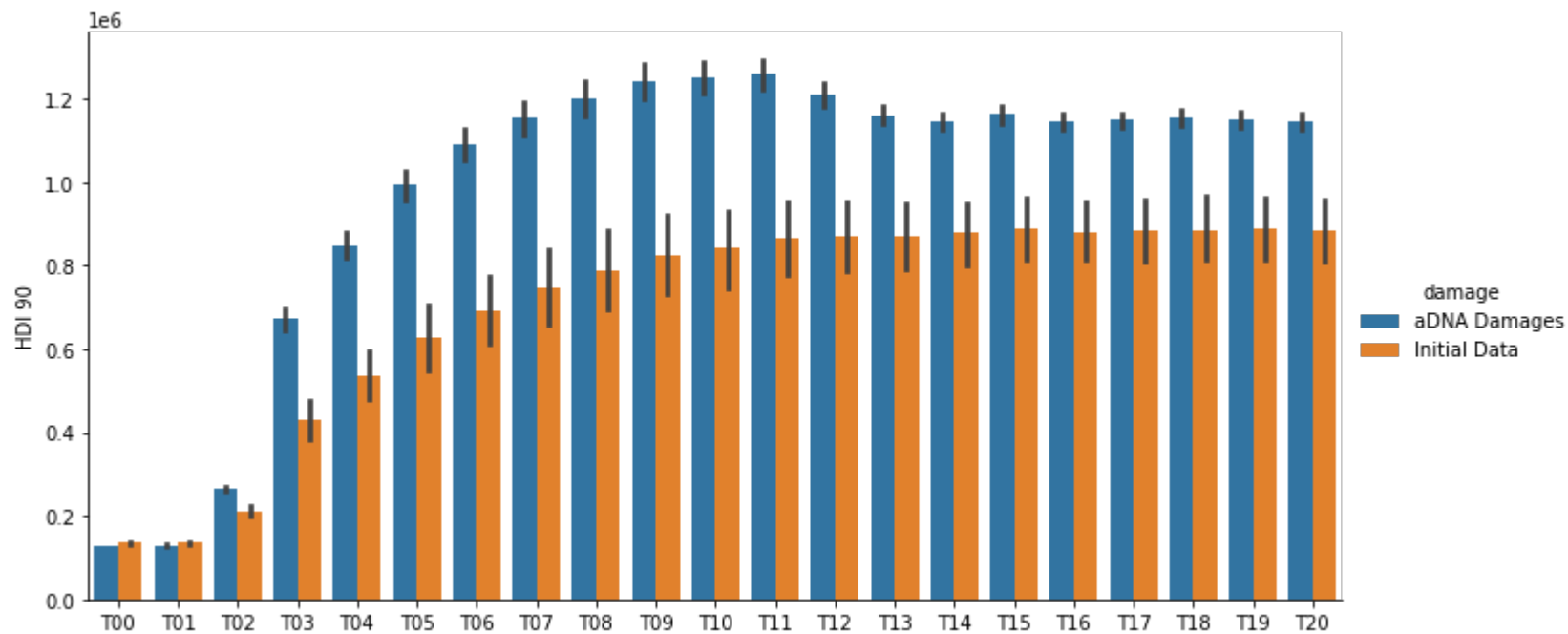# Uncertainty on ancient DNA

# Summary

- Prediction of bacterial population size through time
    - Irrespective of the underlying selection regime and other parameters
- Using `dnadna` package → easy to reuse / reproduce
- Estimation of the aleatoric and epistemic uncertainties

- Transfer learning with aDNA
- Assess interest of transfer learning from other net trained on similar task
- Improve training procedure with SPIDNA (something else than 400 SNP)
- Test on real data

# Thanks

- Flora Jay
- Theophile Sanchez
- Guillaume Charpiat
- Erik M. Bray
- Ben Haller

- Jazeps Medina-Tretmanis

- Maria Avila-Arcos

- Emilia Huerta-Sanchez

- Mathieu Michel