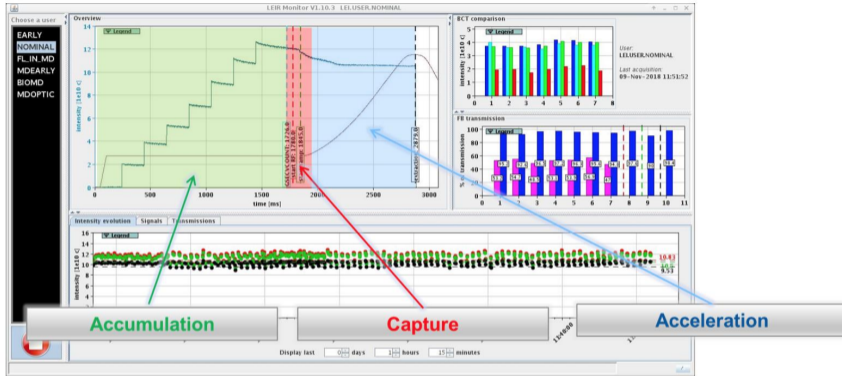


Status of Reinforcement Learning Studies for E-Cooler Operation

N. Biancacci V. Kain A. Latina ■ N. Madysa

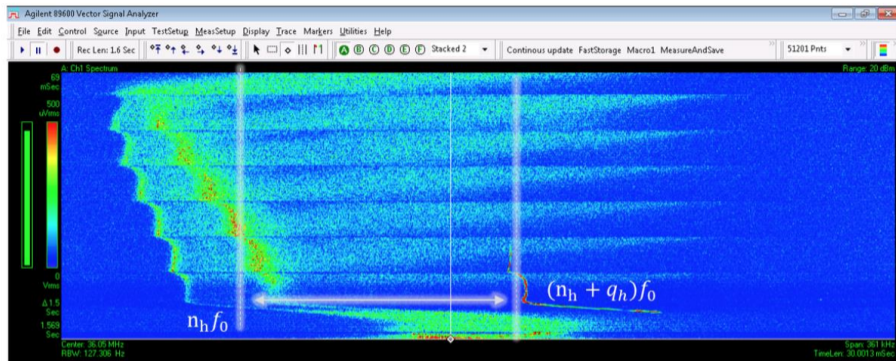
E-BEAM Meeting, 9 September 2020

LEIR Operation



- accumulates injections from Linac 3
- 7 injections every 200 ms, each taking 200 μ s
- cooling throughout flat bottom \Rightarrow fill phase space efficiently
- bunching via RF Capture, acceleration, transfer to PS

LEIR Operation



- accumulates injections from Linac 3
- 7 injections every 200 ms, each taking 200 μ s
- cooling throughout flat bottom \Rightarrow fill phase space efficiently
- bunching via RF Capture, acceleration, transfer to PS

given the previous process:

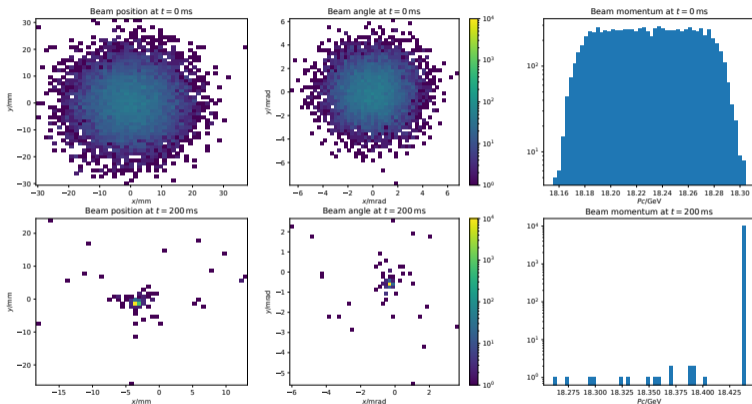
- cooling parameters,
- beam properties before cooling,
- beam properties after cooling,

adjust parameters so that beam after cooling follows some distribution

We will work our way up:

- 1 optimize for minimal $\sigma[p_{\text{out}}]$
- 2 optimize for given $\mu[p_{\text{out}}], \sigma[p_{\text{out}}]$
- 3 optimize other variables than p_{out}
- 4 work with distributions
- 5 work with Schottky spectra

Training Data



- simulation via **RFTTrack**
- simulates the cooler, ignores rest of the accelerator

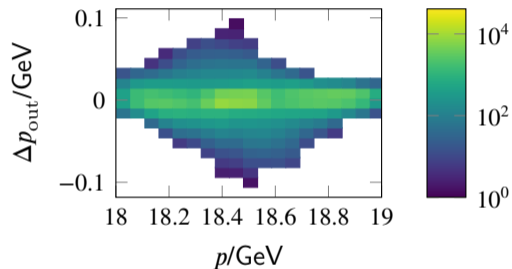
input: random sample of macroparticles

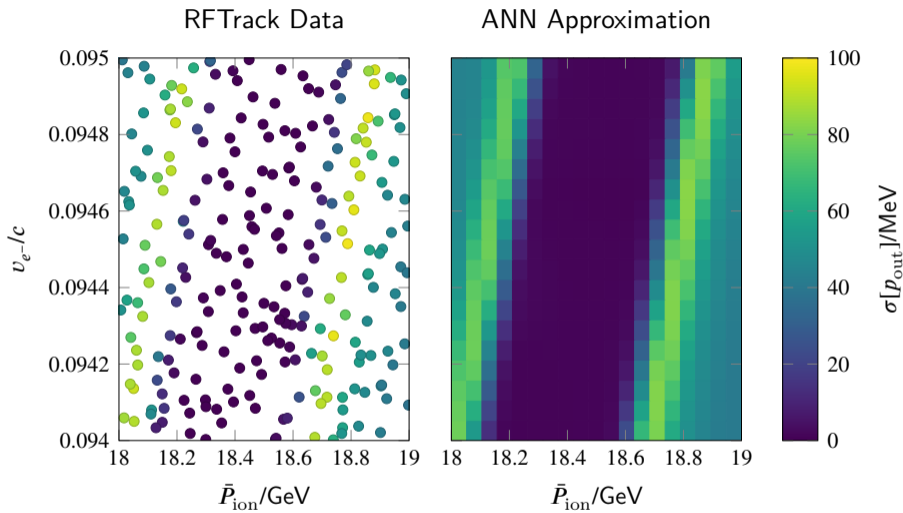
output: cooled macroparticles

Momentum Predictor

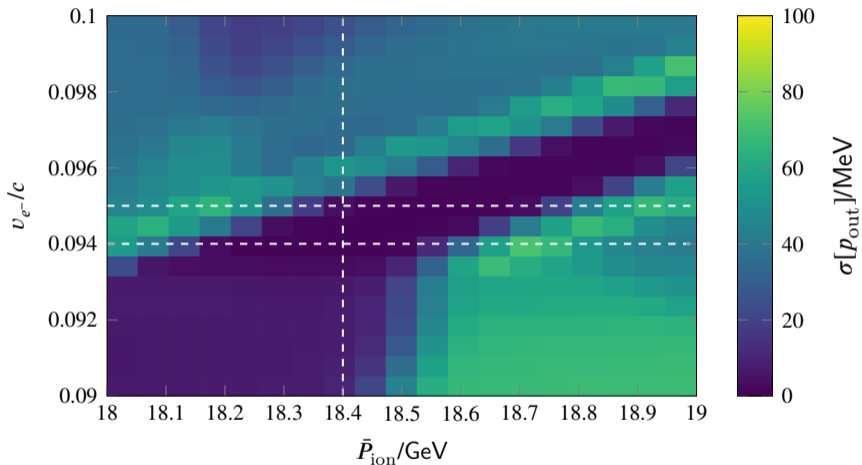
- problem: RFTrack is very complex, takes a long time to complete
- solution: train an ANN to predict outcome of RFTrack simulation

⇒ supervised learning, function approximation

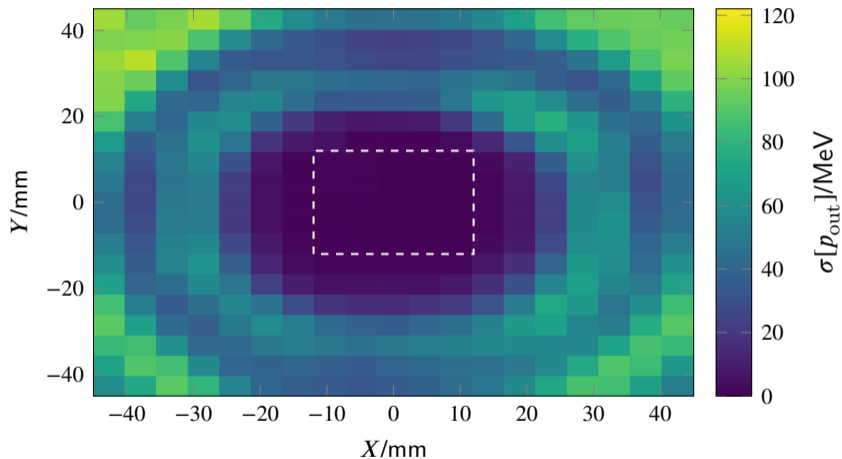


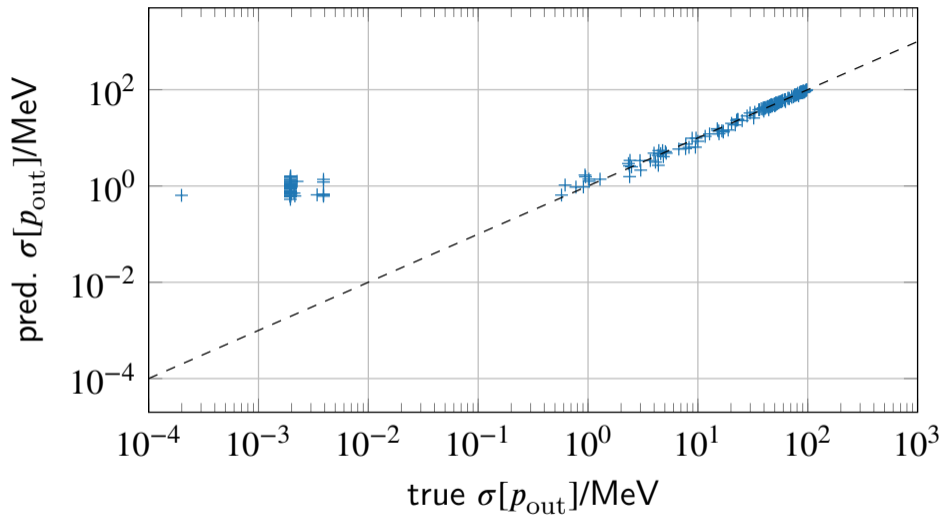


ANN Approximation (extended)



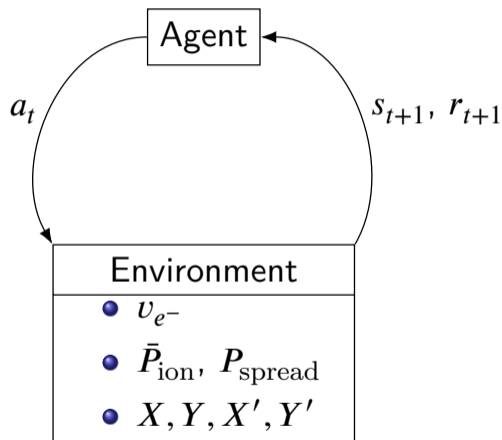
ANN Approximation (extended)





Cooling Agent

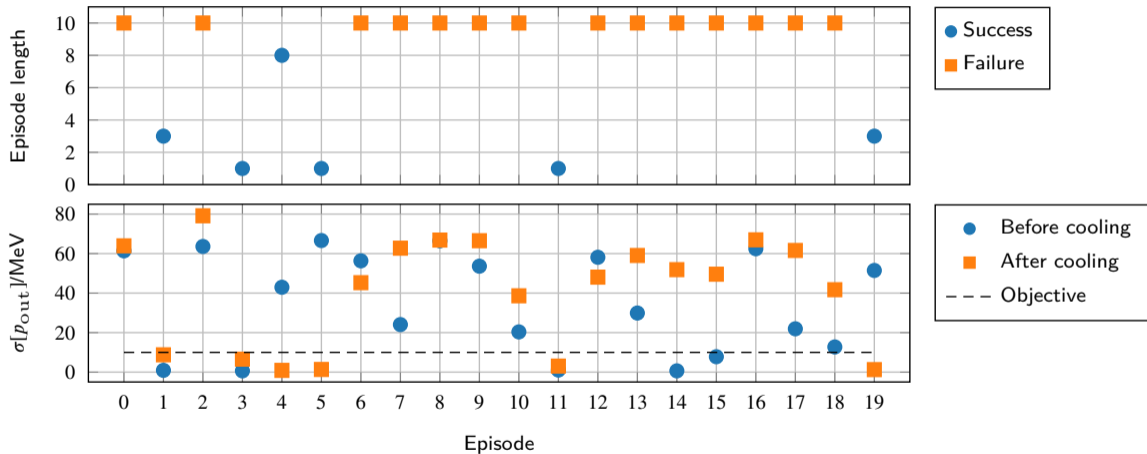
- with predictor, we can generate as many data samples as we need
- use it to train **reinforcement learning** agent
- agent sends cooler settings corrections ($\Delta v_{e^-}, \dots$) to environment
- environment returns next cooling result and a reward



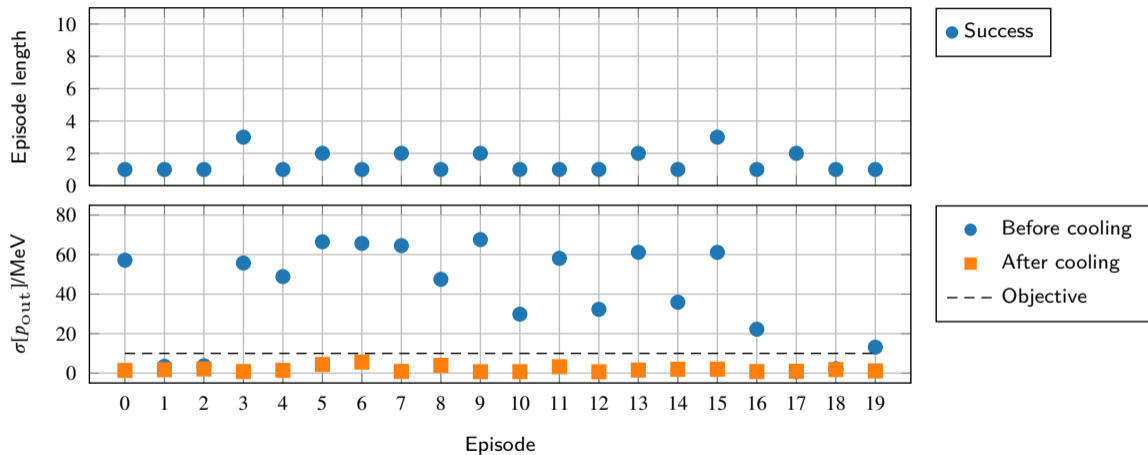
```
1 import gym
2 from stable_baselines3 import TD3
3
4 class LeirEnv(gym.Env): ...
5
6 env = LeirEnv(...)
7 agent = TD3('MlpPolicy', env, ...)
8 agent.learn(10000)
9
10 obs = env.reset()
11 done = False
12 while not done:
13     action, _ = agent.predict(obs)
14     obs, reward, done, info = env.step(action)
15     ...
```



Picking Random Actions

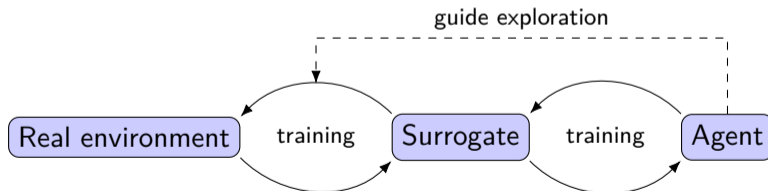


TD3 Algorithm

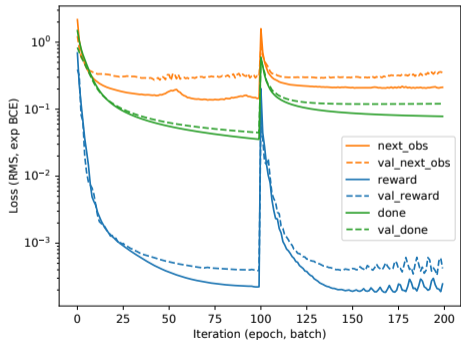


Dyna Architecture

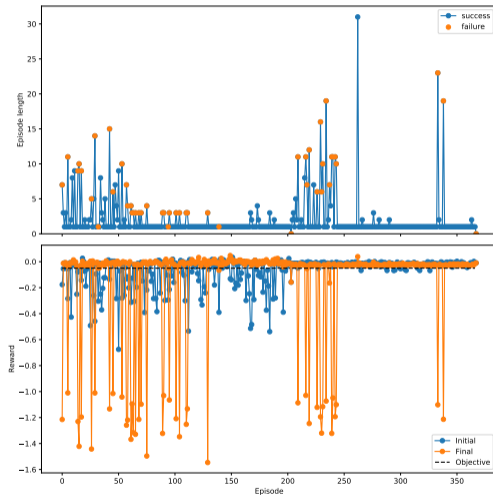
- general approach to improve sample efficiency of model-free agents
- approximate slow real environment with a fast *surrogate* (avoid machine interactions)
- train RL agent on surrogate
- **switch** between training both



Dyna Surrogate Training

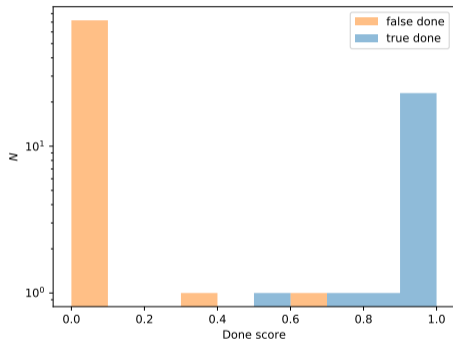


Training of surrogate

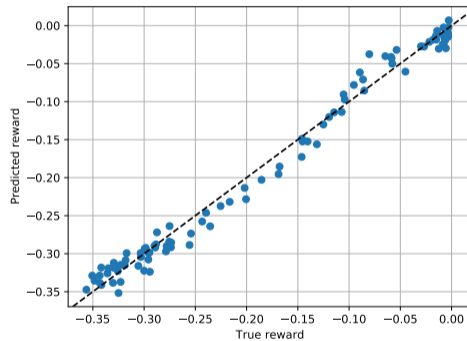


Training of agent on surrogate

Dyna Surrogate Training

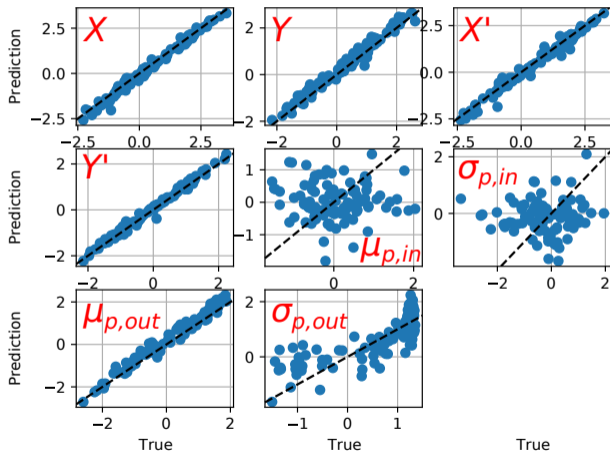


Prediction of end of episode



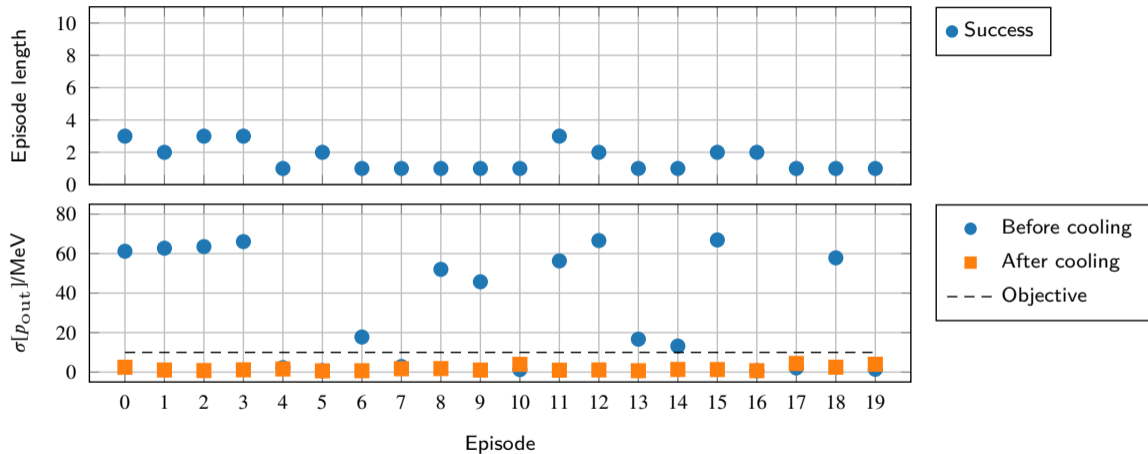
Prediction of reward for each step

Dyna Surrogate Training

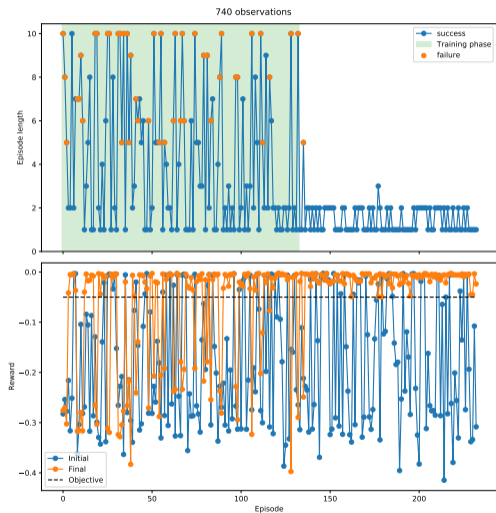


Prediction of next observation (standardized)

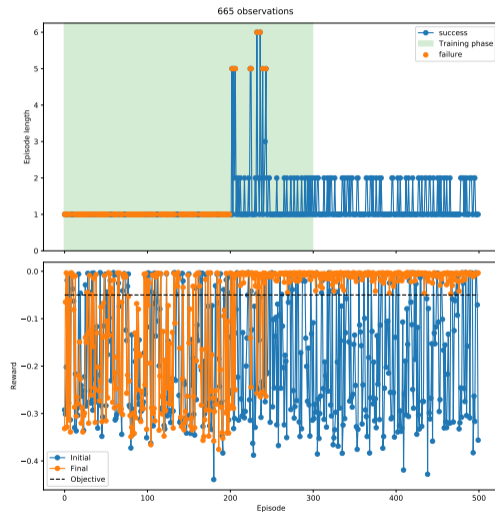
Dyna+TD3



Plain TD3 vs. Dyna+TD3



Plain TD3



Dyna+TD3

Conclusions:

- replaced RFTrack simulation with an ANN
- trained RL agent on simulated data
- improved sample efficiency with model-based RL

Next steps:

- transition predictor from individual particles to beam statistics
- extend optimization problem: $\sigma[x_{\text{out}}]$, $\sigma[y_{\text{out}}]$, $\mu[p_{\text{out}}]$

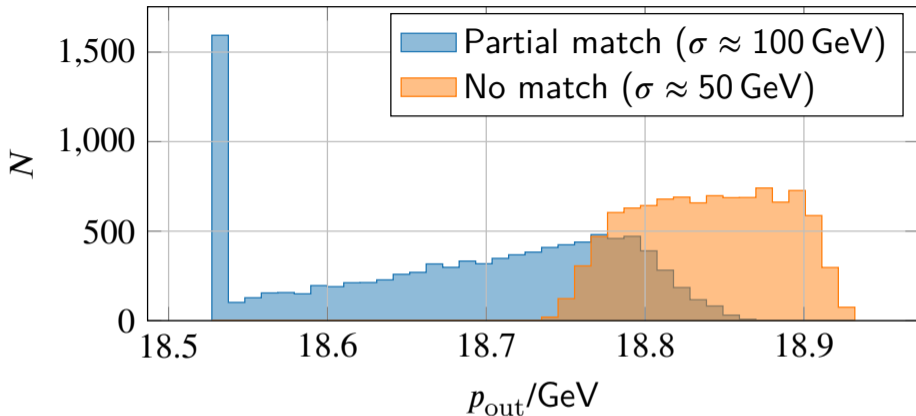
Eventually:

- transition to full distributions of p , x , y , ...
- transition to Schottky spectra



Backup





Reason for ridge structure:

- distribution stays the same if β_{e^-} and \bar{P}_{ion} don't match
- distribution **is torn apart** if β_{e^-} and \bar{P}_{ion} match partially