

Rubén Martínez Cantín

Dpto. Informática e Ingeniería de Sistemas.

(Decision making and RL Break)

- **Bayesian decision making**
 - Active learning
 - Bayesian optimization





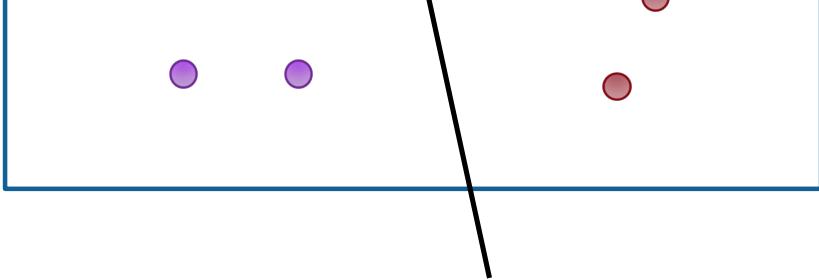
vestido de mujer azul
y negro



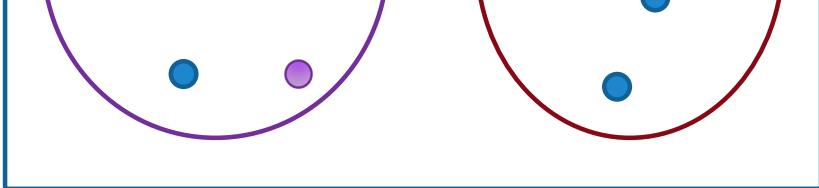
polo de hombre a rayas
azules y negras

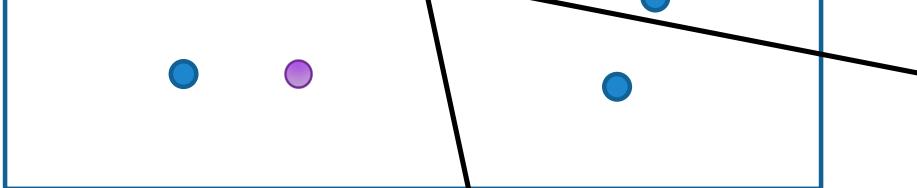


- We need new paradigm of Small Data.









- And we can interact with data

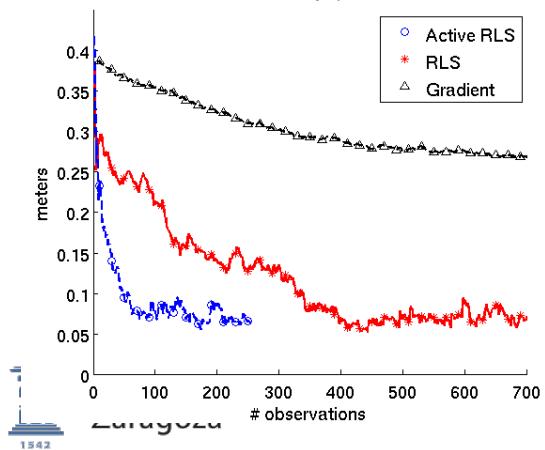
possible data. **Good data is often better than a lot of data.**

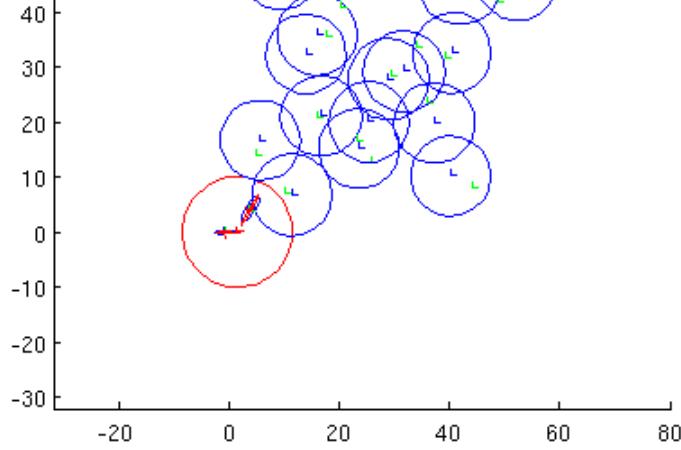


- Humans also query the environment continuously...



Real robot average position error

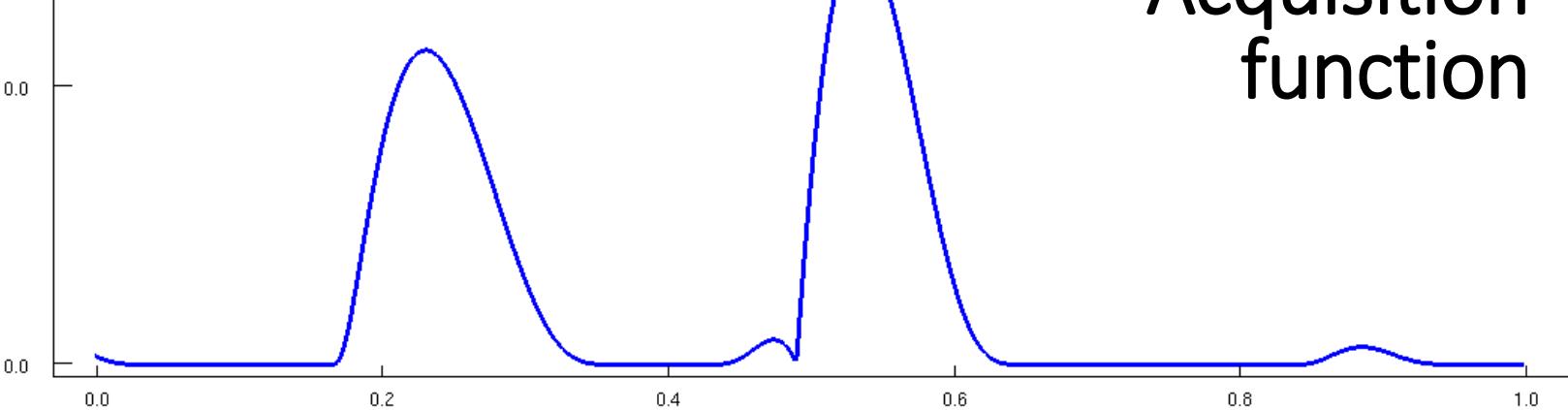


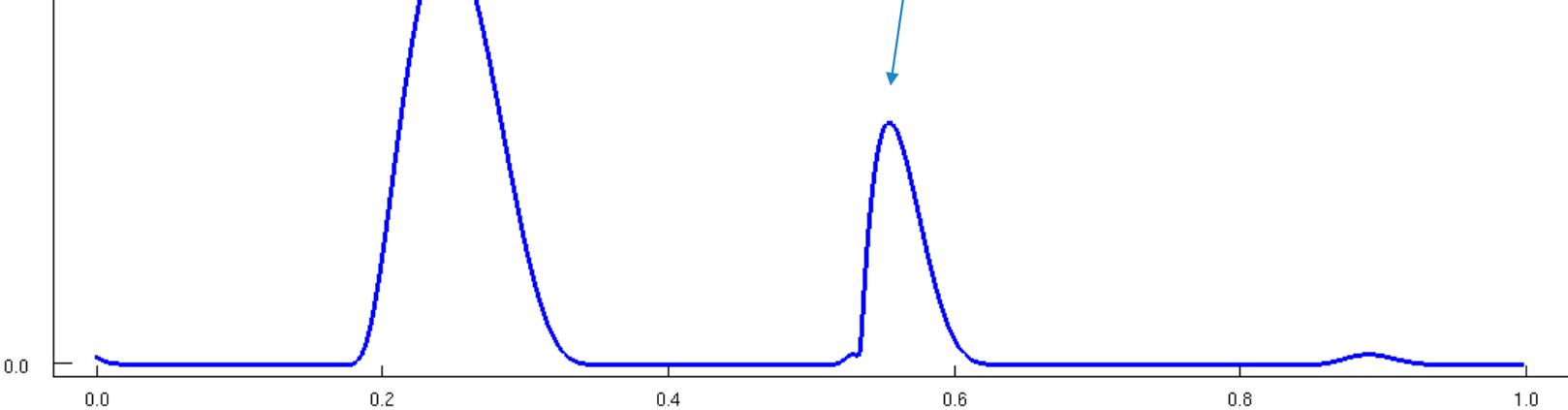


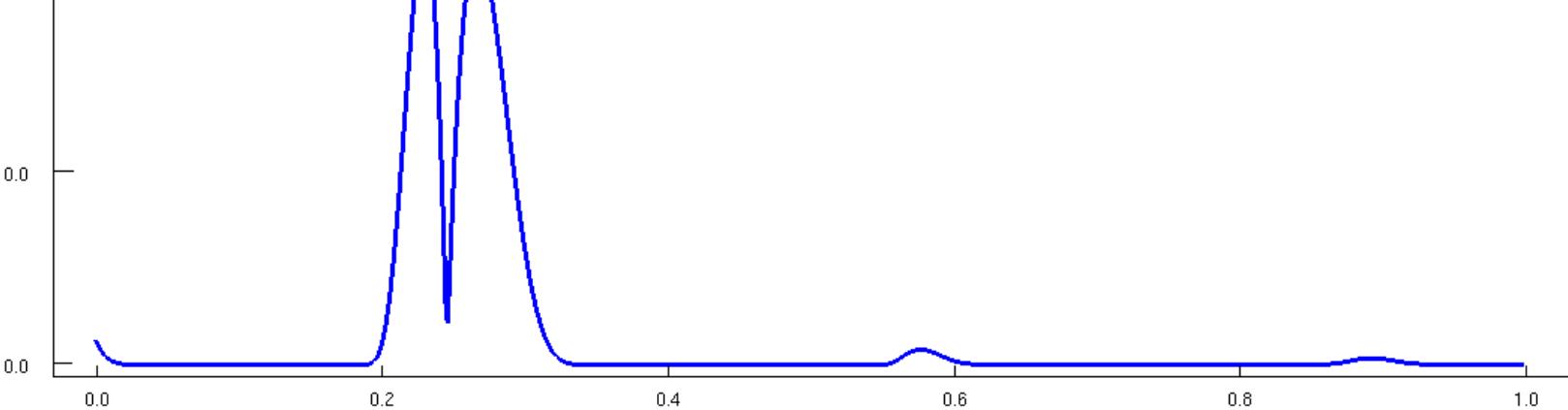
Active policy search: Bayesian optimization



Acquisition function







■ Expected utility > e.g.. Expected Improvement

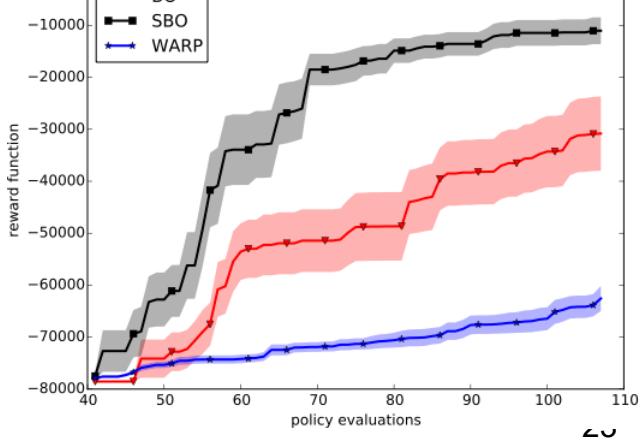
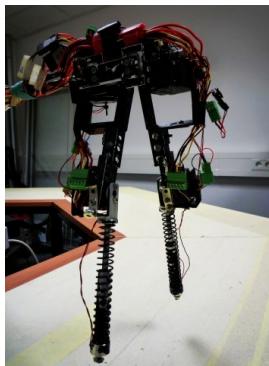
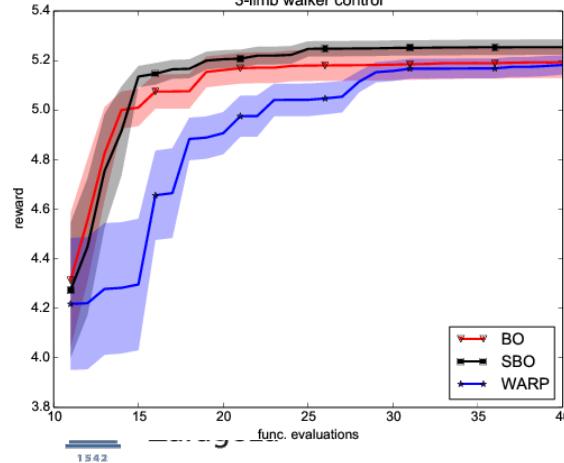
$$\mathbb{E}[I(x)] = \underbrace{(\mu(x) - f(x_{best}))}_{\textit{exploit}} \Phi(d) + \underbrace{\sigma^2(x)}_{\textit{explore}} \phi(d)$$

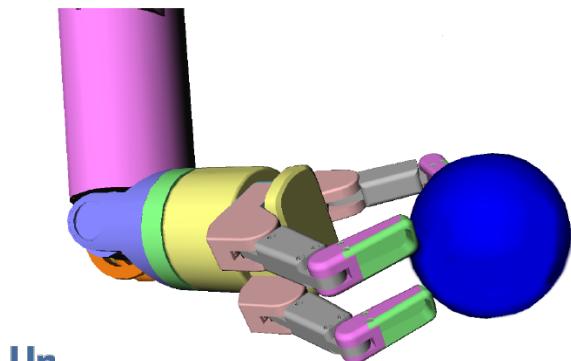
$$d = \frac{x - \mu(x)}{\sigma(x)}$$

$$x_i = \arg\max_{x_i} \phi(x_i; (\mathbf{X}, \mathbf{y})).$$

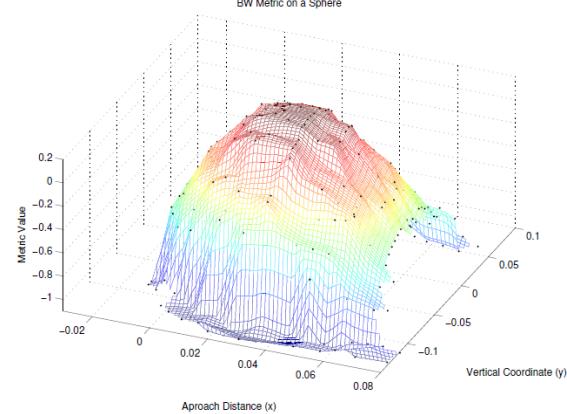
- Augment the data with the new point and response $y_i = f(x_i) + \epsilon$:
 $\mathbf{X} \leftarrow \{\mathbf{X} \cup \mathbf{x}_i\}$ $\mathbf{y} \leftarrow \{\mathbf{y} \cup y_i\}$ $i \leftarrow i + 1$



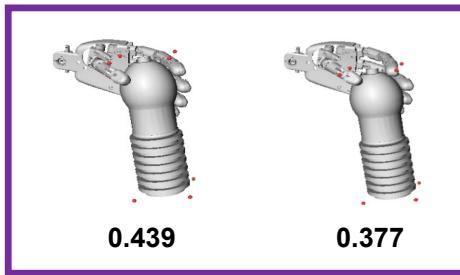
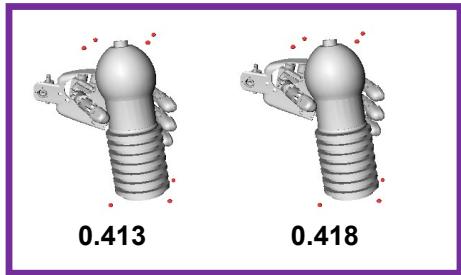


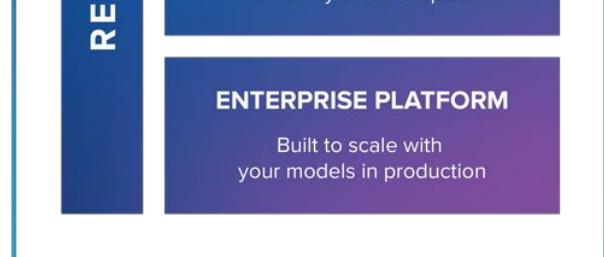
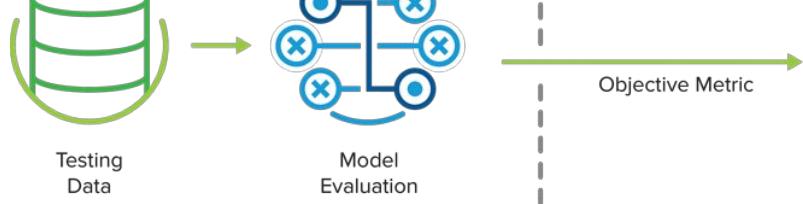


Un
Zaragoza

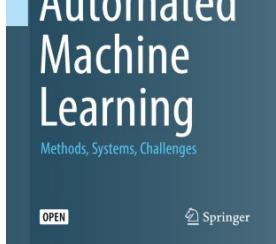


```
[2]: 22:19:42.169855 INFO: [2](1000,1000)  
[2]: 22:19:42.169863 INFO: Latin hypercube sampling
```









Yutian Chen, Aja Huang, Ziyu Wang, Ioannis Antonoglou, Julian Schrittwieser,
David Silver & Nando de Freitas

DeepMind, London, UK
yutianc@google.com

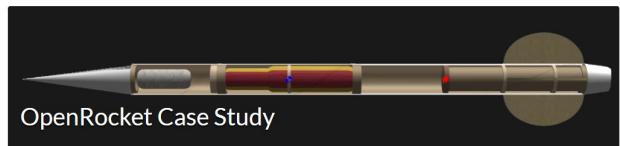
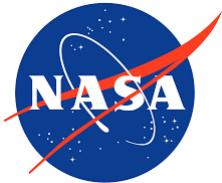


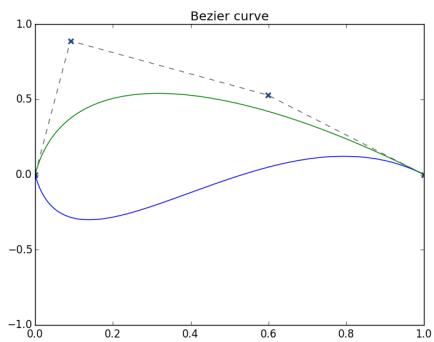
Abstract

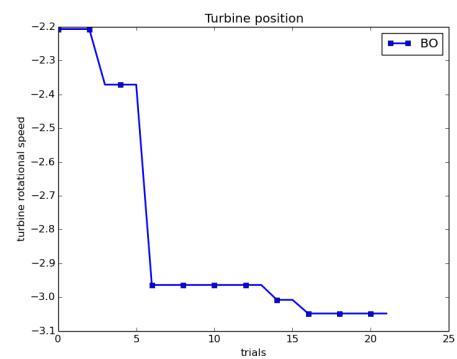
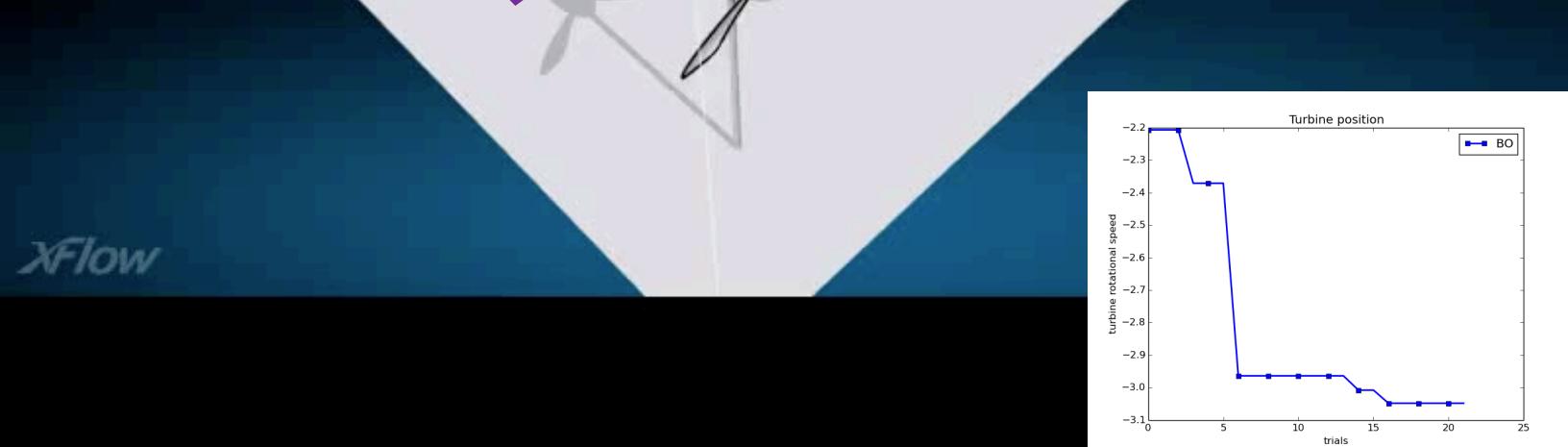
During the development of AlphaGo, its many hyper-parameters were tuned with Bayesian optimization multiple times. This automatic tuning process resulted in substantial improvements in playing strength. For example, prior to the match with Lee Sedol, we tuned the latest AlphaGo agent and this improved its win-rate from 50% to 66.5% in self-play games. This tuned version was deployed in the final match. Of course, since we tuned AlphaGo many times during its development cycle, the compounded contribution was even higher than this percentage. It is our hope that this brief case study will be of interest to Go fans, and also provide Bayesian optimization practitioners with some insights and inspiration.

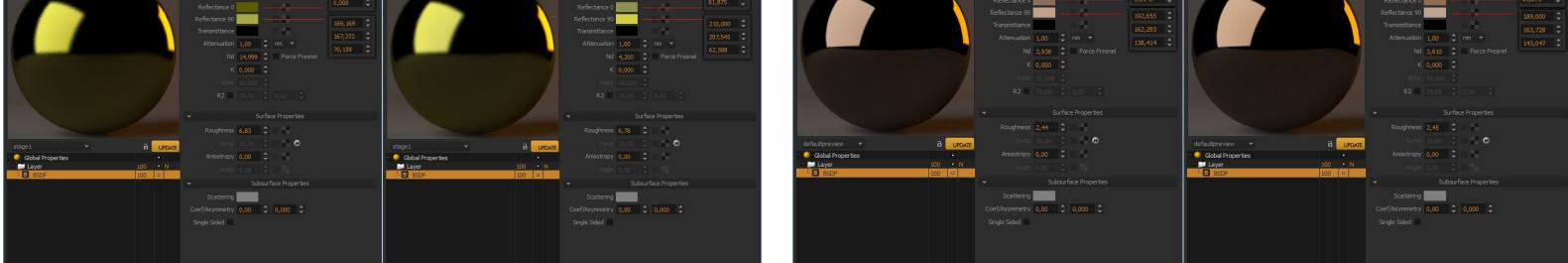


Σ SIGOPT





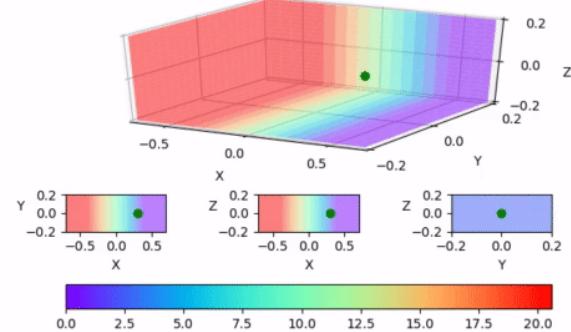
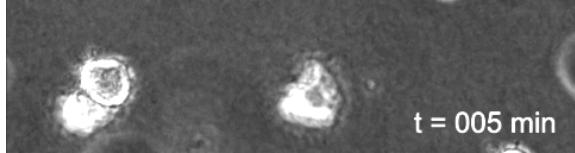


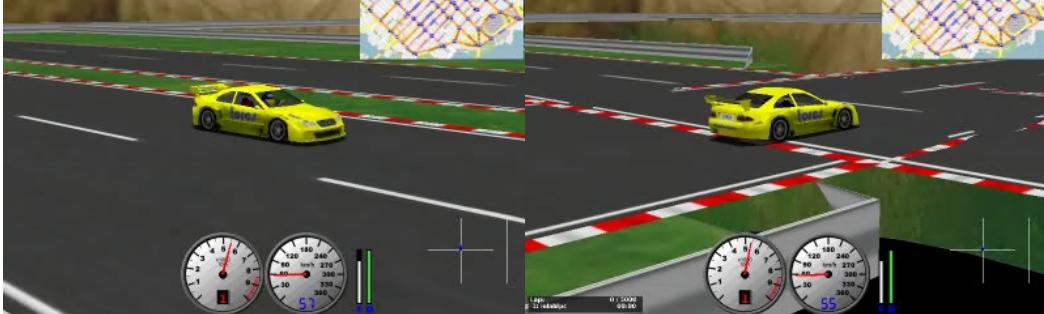


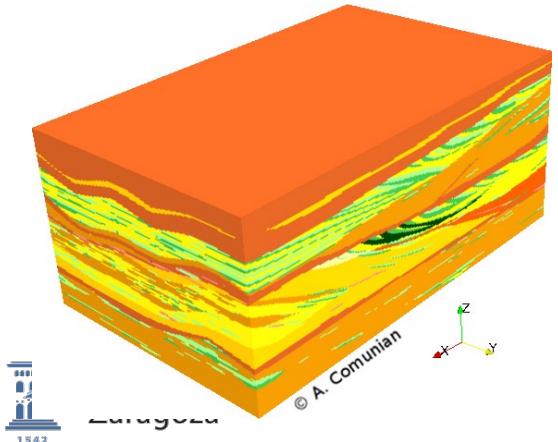
Universidad
Zaragoza

1542

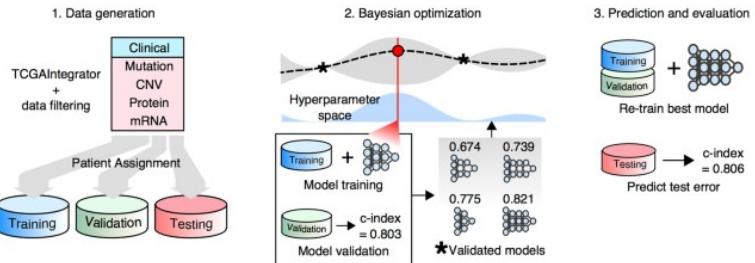
34







C





Zaragoza



50

