



**INNOVACIÓN
INGENIERÍA UAI**
UNIVERSIDAD ADOLFO IBÁÑEZ

FACULTAD DE
INGENIERÍA Y CIENCIAS



SMART +
SUSTAINABLE

Using Quantile Forest For Robust Scheduling Of Astronomic Images Processing:

Inform's Annual Meeting 2022
Indianapolis

Gianfranco Speroni

gsperoni@alumnos.uai.cl

Luis Aburto

luis.aburto@uai.cl

Rodrigo Carrasco

rodrigo.carrascos@uai.cl

October 2022



INNOVACIÓN
INGENIERÍA UAI
UNIVERSIDAD ADOLFO IBÁÑEZ

FACULTAD DE
INGENIERÍA Y CIENCIAS



ALMA OBSERVATORY: THE SCHEDULING PROBLEM TO SOLVE

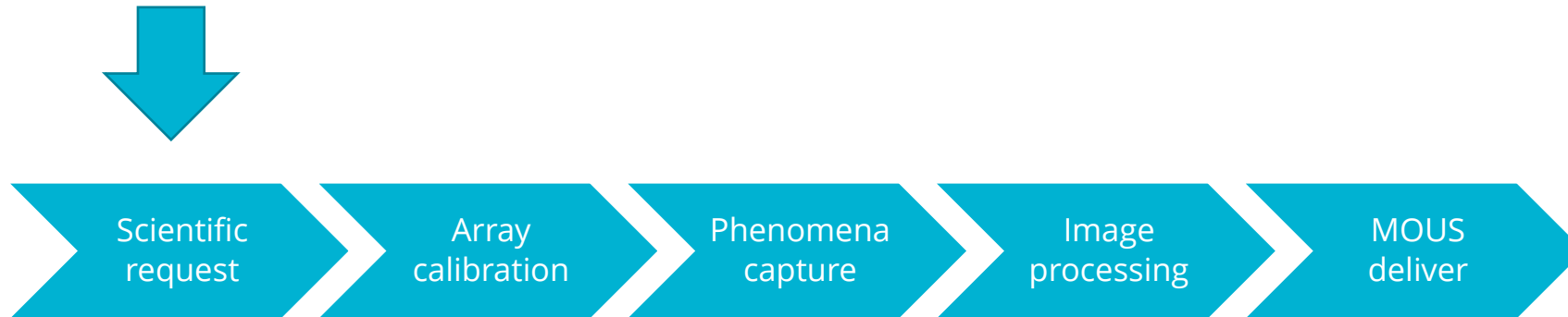
ALMA OBSERVATORY

SOME FACTS

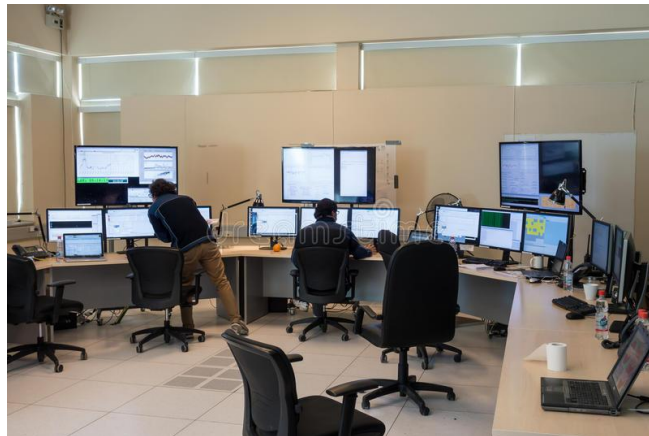
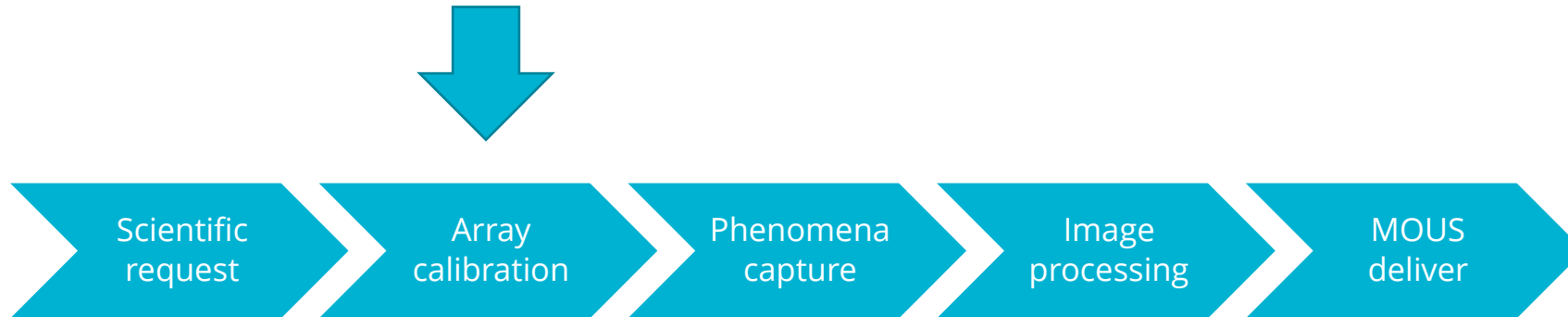
- It's name comes from Atacama Large Millimeter/submillimeter Array (ALMA)
- It's the biggest astronomic observatory in the world (made up of **66** antenas)
- Originated from the colaboration between Europe, North America, Asia and the Chilean Republic
- It works 24/7 all 365 days of the year
- Most of the work done in the ALMA installations is related with: Stars formation, molecular clouds and the early universe



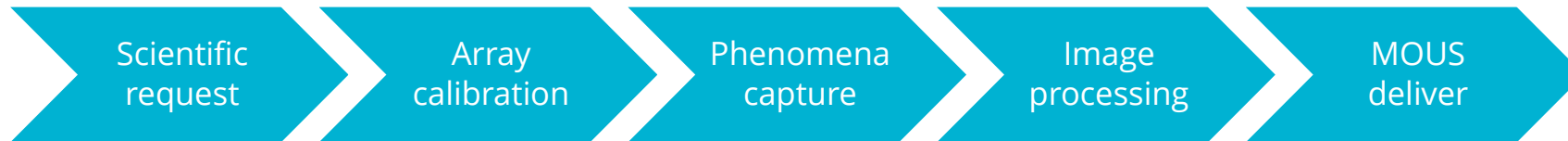
BUSINESS PROCESS



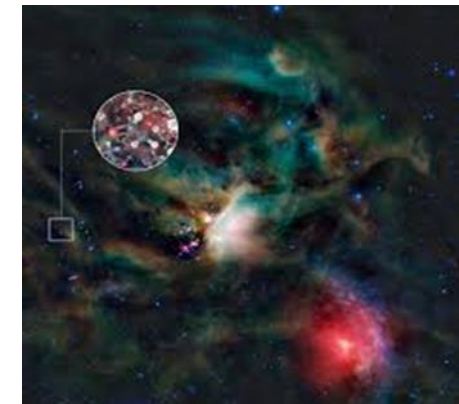
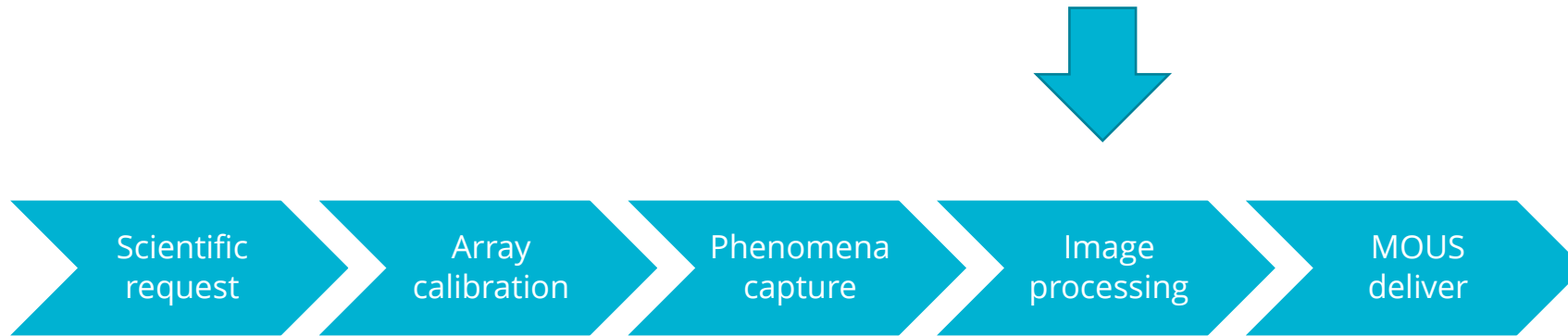
BUSINESS PROCESS



BUSINESS PROCESS

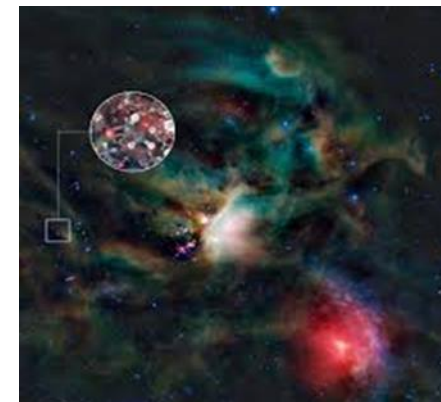
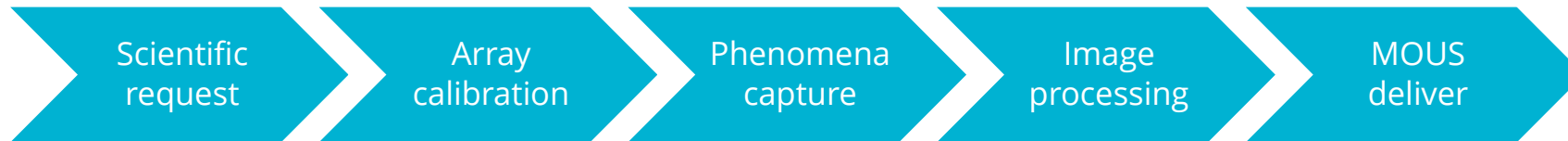


BUSINESS PROCESS

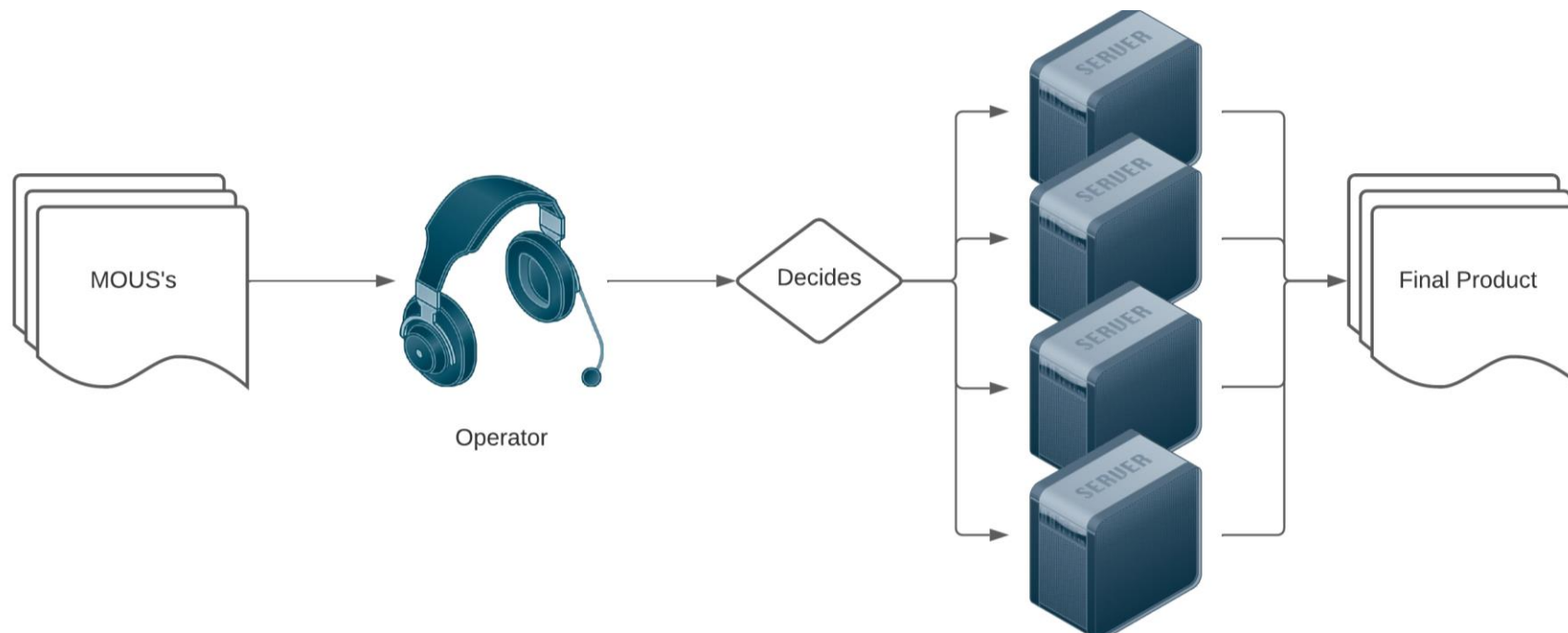


BUSINESS PROCESS

FACULTAD DE
INGENIERÍA Y CIENCIAS



BUSINESS PROBLEM



→ALMA gives a warranty of 30 days to deliver the final product of the scientific request

→The operator doesn't know exactly how much the image processing of each MOUS it's going to last

→The system has a fixed level of service (number of available servers)

→Therefore finding an efficient way to assign and sequenciate the released jobs to the servers is key!

LITERATURE REVIEW

UNRELATED MACHINES SCHEDULING WITH STOCHASTIC PROCESSING TIMES (SKUTELLA, 2016)

Presents a MIP for completion time minimization and also a strategy to derive in a LP-relaxation

FROM PREDICTIVE TO PRESCRIPTIVE ANALYTICS (BERTSIMAS, 2020)

Proposes an approach to unify an ML model and a optimization algorithm with a weighting function that is meant to reduce the performance differences

OBJECTIVES

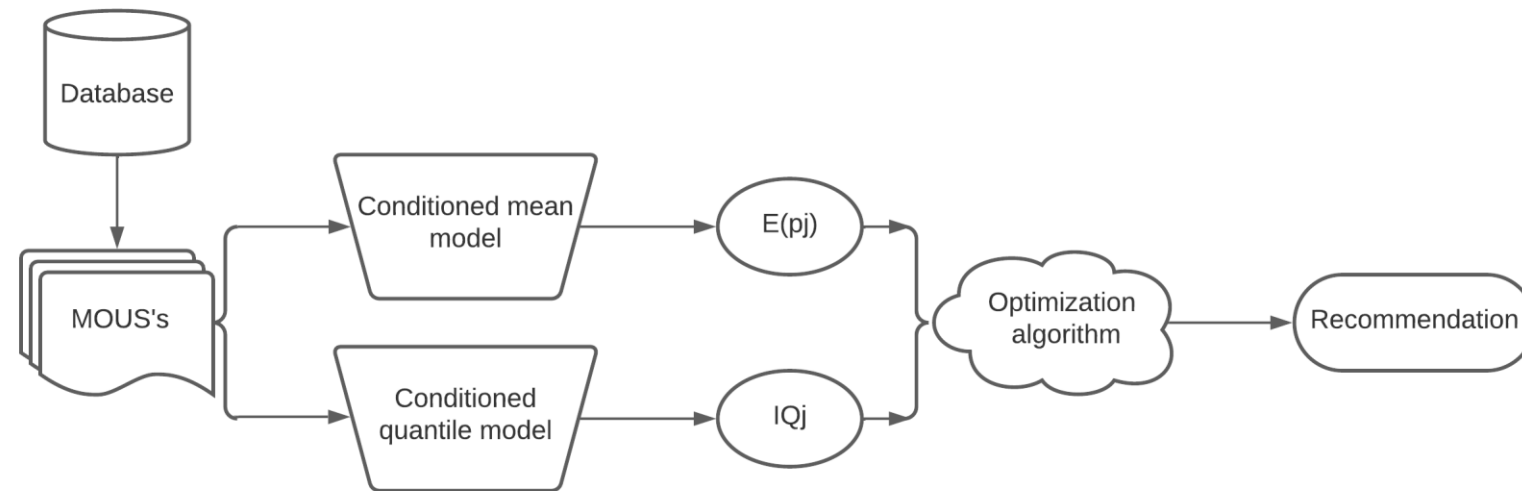
GENERAL OBJECTIVE

Minimize the flow times of the MOUS (member of unit set) and understand which are the key factors to estimate the processing times and their uncertainty

SPECIFIC OBJECTIVES

- Generate estimations for the processing times of the MOUS
- Generate estimations over the uncertainty related to the estimations of the processing time of the MOUS
- Understand which variables explain better the estimation of processing time of the MOUS and which variables explain better the uncertainty related to the estimations of the processing time of the MOUS
- Formulate a stochastic and robust mathematical model of assignation and sequenciation capable of getting solutions that are near in relation to the optimal solution, by considering the presence of the uncertainty related to the estimations

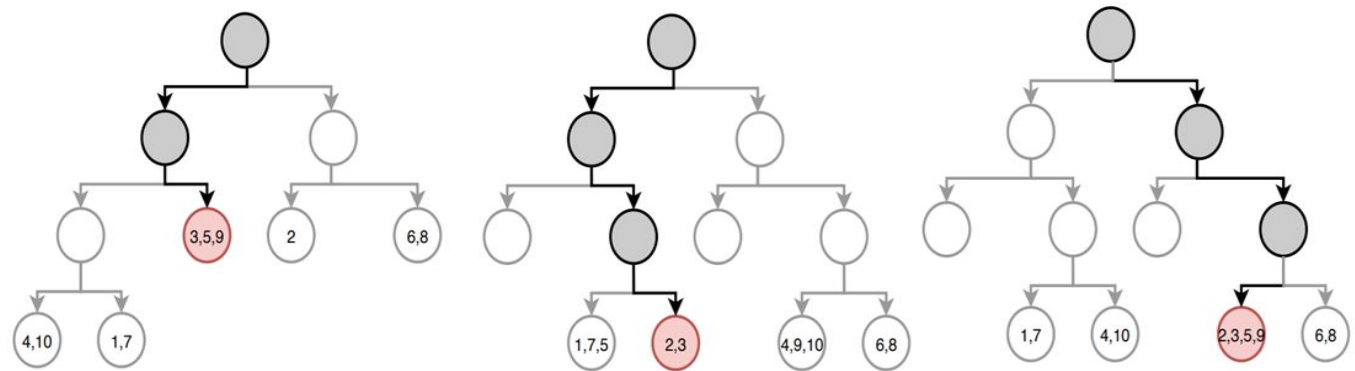
METHODOLOGY



PREDICTIVE MODEL

QUANTILE FOREST

- The elaboration of the quantile forest is no different than the elaboration of random forest. The main difference relies in the moment when a prediction is made
- Given 'k' bagged trees (like in random forest) we can compute a quantile prediction for some observation 'x' by considering the aggregation of each tree estimation (all 'k' estimations) for the unknown Y as an empirical distribution function for Y
- Then for a specified percentile we do $Q_{(Y|X)}(\tau) = \inf\{y: F_{(Y|X)}(y) \geq \tau\}$



OPTIMIZATION MODELS

ELEMENTS OF THE FORMULATION:

- M : Amount of machines
- J : Amount of jobs
- T : Time horizon
- τ : Window of time to assign the start of the j job processing
- x_{ijt} : Assignment variable that indicates the start of the j job processing at the machine i in the t moment
- p_j : Processing time of the j job
- r_j : Release time of the j job

PERFECT INFORMATION FORMULATION:

$$\text{Min} \sum_{j \in J} F_j$$

$$\sum_{i \in M} \sum_{t \in T} x_{ijt} = 1 \quad \forall j \in J \quad (1)$$

$$\sum_{j \in J} \sum_{\tau = \max(0, t - p_j)}^{t-1} x_{ijt} \leq 1 \quad \forall i \in M, t \in T \quad (2)$$

$$F_j = \sum_{i \in M} \sum_{t \in T} x_{ijt} (t + p_j - r_j) \quad \forall j \in J \quad (3)$$

$$\sum_{i \in M} \sum_{t=0}^{r_j-1} x_{ijt} = 0 \quad \forall j \in J \quad (4)$$

$$x_{ijt} \in \{0,1\} \quad \forall i \in M, j \in J, t \in T \quad (5)$$

OPTIMIZATION MODELS

ELEMENTS OF THE FORMULATION:

- M : Amount of machines
- J : Amount of jobs
- T : Time horizon
- τ : Window of time to assign the start of the j job processing
- x_{ijt} : Assignment variable that indicates the start of the j job processing at the machine i in the t moment
- \hat{p}_j : Processing time of the j job
- r_j : Release time of the j job
- \widehat{IQ}_j : Interquantile range for the processing time of the j job
- λ : Penalization level

RECOMMENDATION FORMULATION WITH ABSOLUTE UNCERTAINTY PENALIZATION:

$$\text{Min} \sum_{j \in J} F_j + \lambda * \widehat{IQ}_j$$

$$\sum_{i \in M} \sum_{t \in T} x_{ijt} = 1 \quad \forall j \in J \quad (1)$$

$$\sum_{j \in J} \sum_{\tau = \max(0, t - \hat{p}_j)}^{t-1} x_{ijt} \leq 1 \quad \forall i \in M, t \in T \quad (2)$$

$$F_j = \sum_{i \in M} \sum_{t \in T} x_{ijt} (t + \hat{p}_j - r_j) \quad \forall j \in J \quad (3)$$

$$\sum_{i \in M} \sum_{t=0}^{r_j-1} x_{ijt} = 0 \quad \forall j \in J \quad (4)$$

$$x_{ijt} \in \{0,1\} \quad \forall i \in M, j \in J, t \in T \quad (5)$$

OPTIMIZATION MODELS

ELEMENTS OF THE FORMULATION:

- M : Amount of machines
- J : Amount of jobs
- T : Time horizon
- τ : Window of time to assign the start of the j job processing
- x_{ijt} : Assignment variable that indicates the start of the j job processing at the machine i in the t moment
- \hat{p}_j : Processing time of the j job
- r_j : Release time of the j job
- \widehat{IQ}_j : Interquantile range for the processing time of the j job
- λ : Penalization level

RECOMMENDATION FORMULATION WITH EXPONENTIAL UNCERTAINTY PENALIZATION:

$$\text{Min} \sum_{j \in J} F_j + \widehat{IQ}_j^\lambda$$

$$\sum_{i \in M} \sum_{t \in T} x_{ijt} = 1 \quad \forall j \in J \quad (1)$$

$$\sum_{j \in J} \sum_{\tau = \max(0, t - \hat{p}_j)}^{t-1} x_{ijt} \leq 1 \quad \forall i \in M, t \in T \quad (2)$$

$$F_j = \sum_{i \in M} \sum_{t \in T} x_{ijt} (t + \hat{p}_j - r_j) \quad \forall j \in J \quad (3)$$

$$\sum_{i \in M} \sum_{t=0}^{r_j-1} x_{ijt} = 0 \quad \forall j \in J \quad (4)$$

$$x_{ijt} \in \{0,1\} \quad \forall i \in M, j \in J, t \in T \quad (5)$$

OPTIMIZATION MODELS

ELEMENTS OF THE FORMULATION:

- M : Amount of machines
- J : Amount of jobs
- T : Time horizon
- τ : Window of time to assign the start of the j job processing
- x_{ijt} : Assignment variable that indicates the start of the j job processing at the machine i in the t moment
- \hat{p}_j : Processing time of the j job
- r_j : Release time of the j job
- \widehat{IQ}_j : Interquantile range for the processing time of the j job
- λ : Penalization level

RECOMMENDATION FORMULATION WITH PERCENTUAL UNCERTAINTY PENALIZATION:

$$\text{Min} \sum_{j \in J} F_j + \lambda * \frac{\widehat{IQ}_j}{\hat{p}_j}$$

$$\sum_{i \in M} \sum_{t \in T} x_{ijt} = 1 \quad \forall j \in J \quad (1)$$

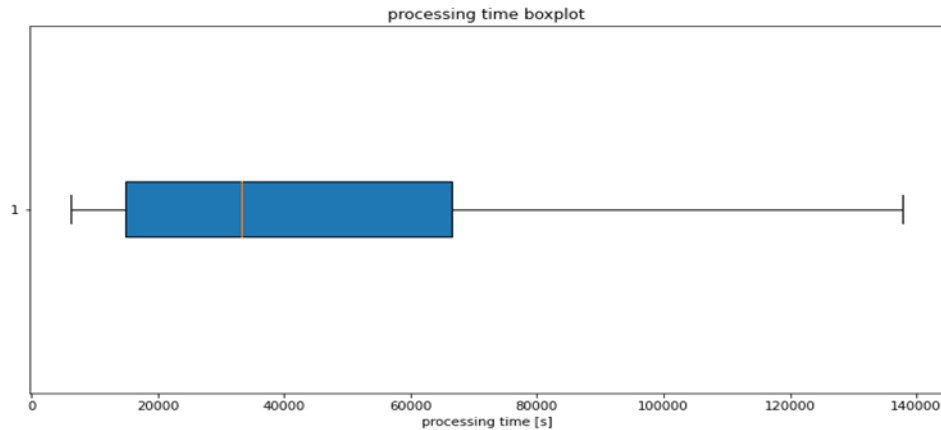
$$\sum_{j \in J} \sum_{\tau = \max(0, t - \hat{p}_j)}^{t-1} x_{ijt} \leq 1 \quad \forall i \in M, t \in T \quad (2)$$

$$F_j = \sum_{i \in M} \sum_{t \in T} x_{ijt} (t + \hat{p}_j - r_j) \quad \forall j \in J \quad (3)$$

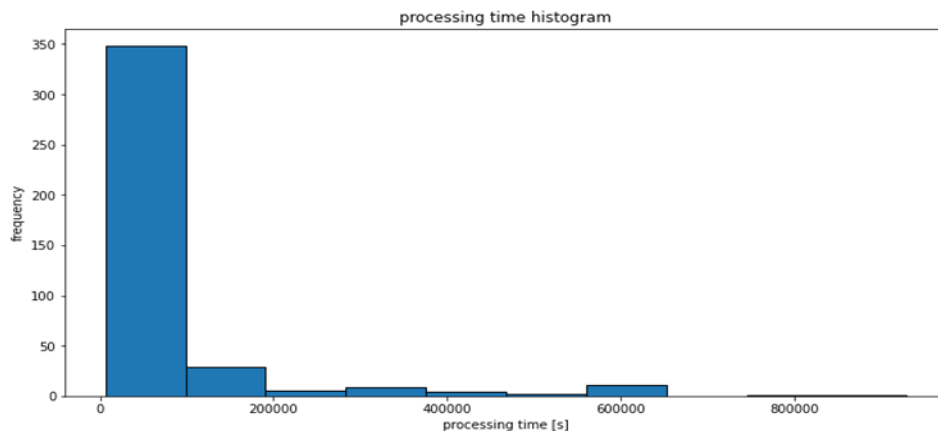
$$\sum_{i \in M} \sum_{t=0}^{r_j-1} x_{ijt} = 0 \quad \forall j \in J \quad (4)$$

$$x_{ijt} \in \{0,1\} \quad \forall i \in M, j \in J, t \in T \quad (5)$$

EMPIRICAL APLICATION TO ALMA SCHEDULING



estadísticos	processing time [h]
count	0,1
mean	20,6
std	35,9
min	1,7
25%	4,1
50%	9,2
75%	18,5
max	258,1



- The support of the processing time is defined in a pretty wide range
- In order to recommend a schedule it's necessary to know how long is gonna take a job to process

DATASET DESCRIPTION

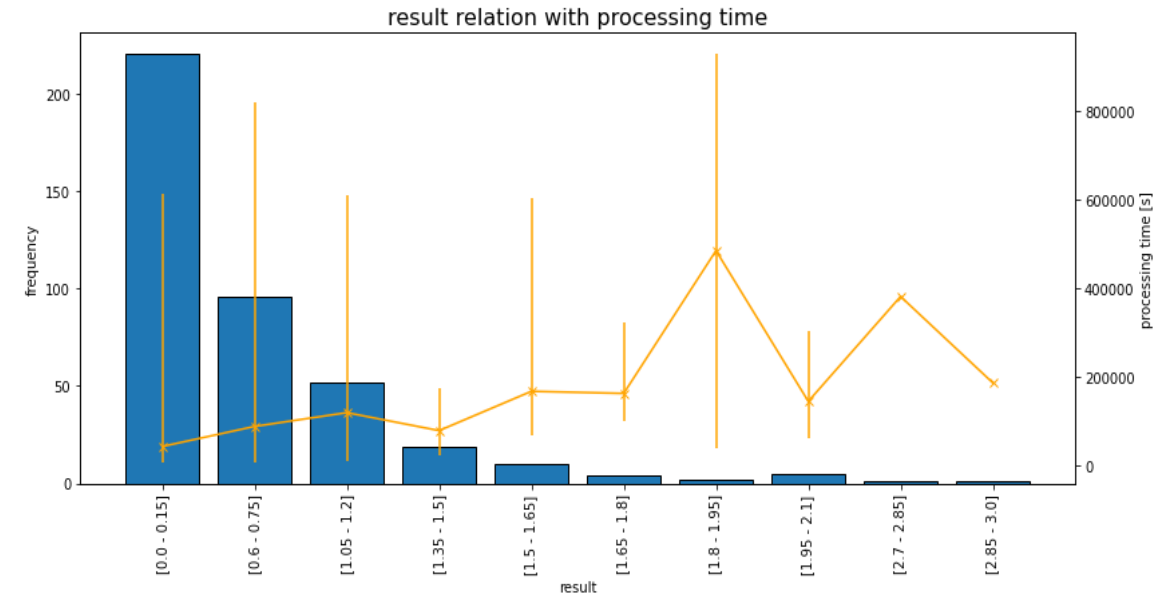
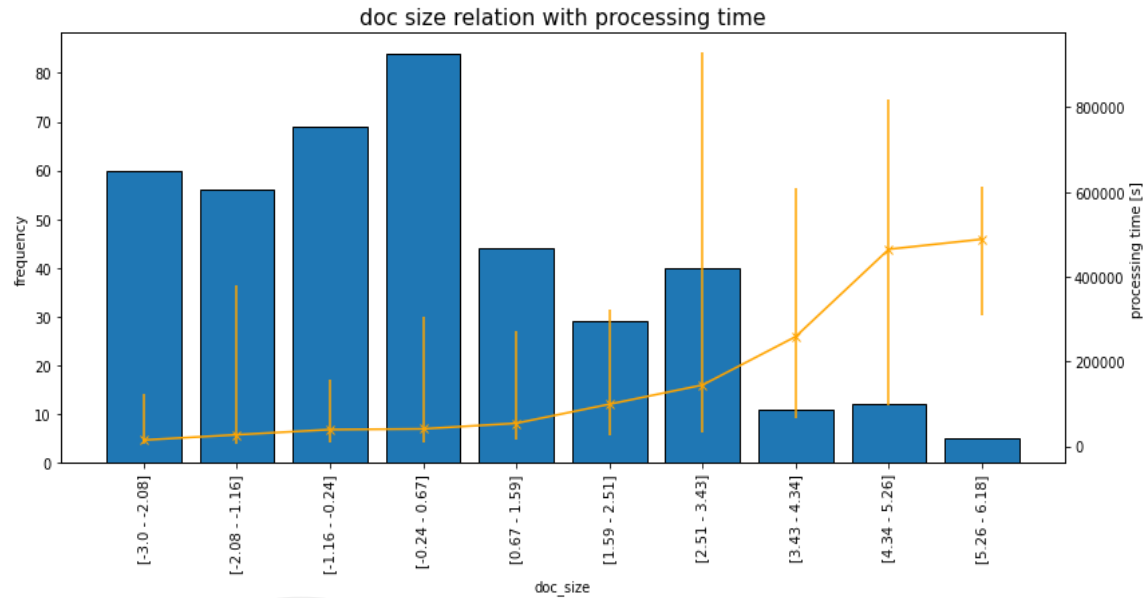
DIMENSIONS

- 16 Attributes
- 410 MOUS

ATTRIBUTES DICTIONARY

- Requested_array: Tells if the array has been requested
- Scheduling_blocks: Is the amount of specific storage units
- Receiver_band: Tells which band is being used
- Antennas: Is the amount of antennas that are being used
- Spectral_windows: Is the amount of defined quadrants
- Fields: The amount of fields that are being used
- Channels: The amount of channels that are being used
- Result: Resume of expert knowledge
- Observation: Time of phenomena observation in seconds
- Doc_size: Size in GB of the MOUS
- Bandpass: Calibration parameter
- Phase: Calibration parameter
- Target: Calibration parameter
- Pointing: Calibration parameter
- Atmosphere: Resume of atmospheric conditions
- Processing_time: Time that a MOUS requires to be processed in the servers

RELATION BETWEEN THE PROCESSING TIME AND THE FEATURES



There is useful information on the MOUS metadata to estimate the processing time !

SELECTING A MODEL TO PREDICT THE PROCESSING TIME

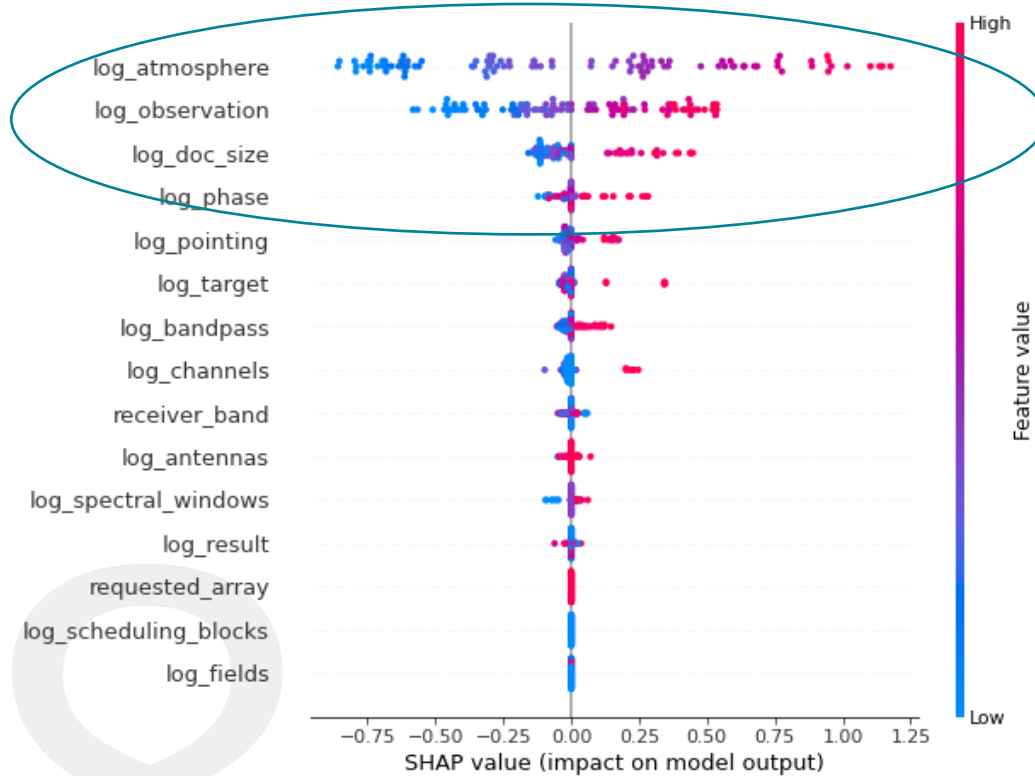
Modelo	Tratamiento	WMAPE (test)	MAPE (test)	R2 (test)	WMAPE (training)	MAPE (training)	R2 (training)
Elastic Net	Standard scaler	38%	40%	72%	36%	39%	77%
K-Neighbors	Standard scaler	28%	22%	71%	1%	0%	100%
SVR	Standard scaler + log en features	24%	19%	77%	8%	7%	95%
Random Forest	NA	29%	21%	73%	11%	8%	96%
Light Gradient Boosting	Standard scaler + log en target + log en features	24%	16%	76%	11%	6%	90%
Gradient Boosting	NA	28%	21%	73%	4%	8%	100%
Adaptative Boosting	Standard scaler	24%	17%	74%	1%	2%	100%

SELECTING A MODEL TO ESTIMATE THE CONDITIONED IQ

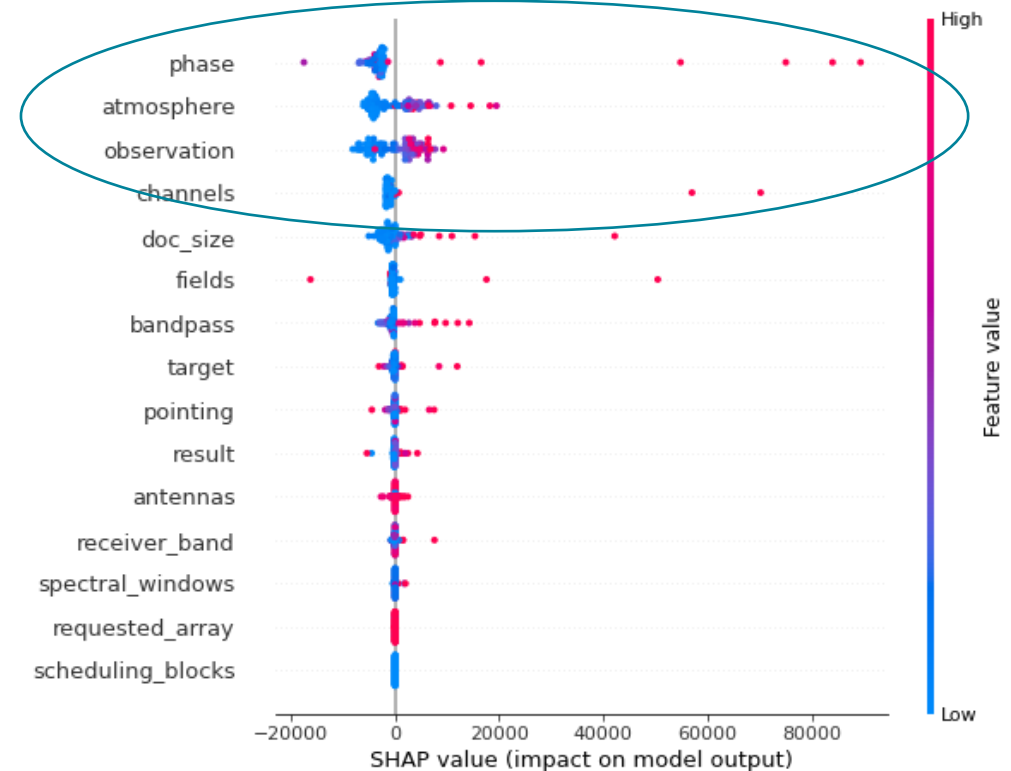
	Quantile forest	Ada Boost Bootstrap
Percentage of data points inside	45%	28%
Mean absolute inter quantile distance	23586	12994
Mean percentual inter quantile distance	23%	11%
Quantile loss 25%	6189	7235
Quantile loss 75%	16474	9835

RANKING THE POWER OF EXPLAINABILITY OF THE VARIABLES FOR THE REGRESSION MODELS

PROCESSING TIME ESTIMATION

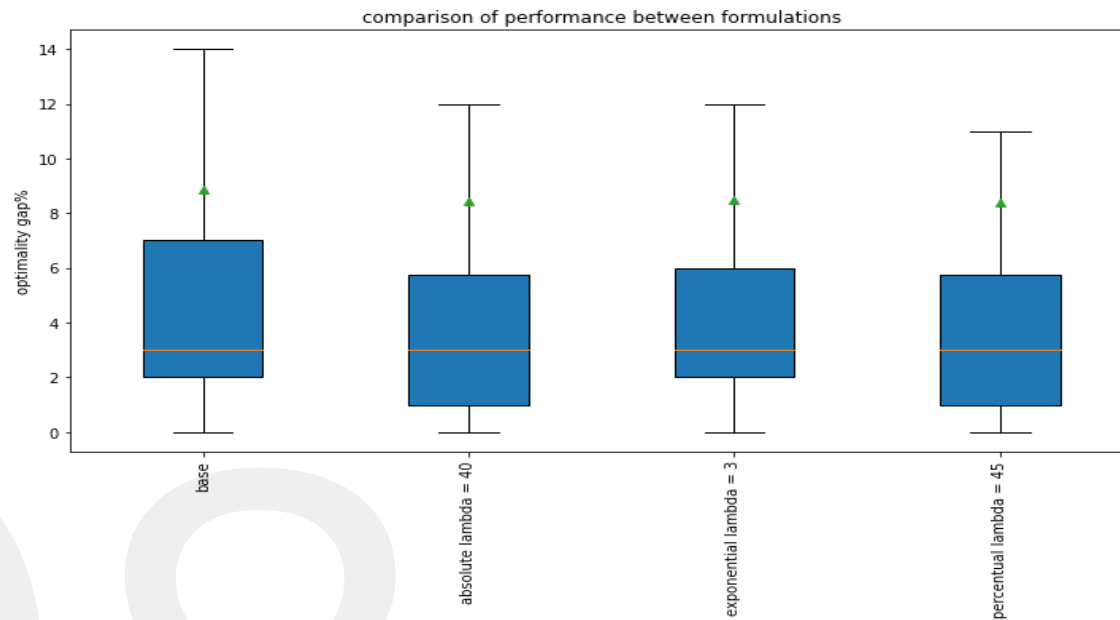


UNCERTAINTY ESTIMATION



The variables that explain the mean of the processing time are different of those which explain the uncertainty related to that estimation!

SELECTING THE BEST OPTIMIZATION ALGORITHM



Compared formulations	sensitibity	mean_gap%	25_gap%	median_gap%	75_gap%
Oracle-Recommendation (Alone)	NA	8,88	2	3	7
Oracle-Recommendation with percentual penalization	45	8,4	1	3	5,75

Empirical improve of
5,4%! (regret)

FINAL COMMENTS

- The conditioned mean model improved its performance by applying nonlinear transformations to both features and response variable
- The variables that explain the expectation of the variable of interest differ from the ones that explain the uncertainty related to the expectation
- Including the uncertainty in the optimization algorithm reported a 5.4% mean improve in the optimality gap

FUTURE WORK

PENALIZING THE UNCERTAINTY ON THE CONSTRAINTS OF THE MIP

In order to have an effect on the decision variable, one approach is to apply the same penalization that has been exposed and place it on the constraints that contained \hat{p}_j

CONSTRAINT FOR UNCERTAINTY BALANCING ACROSS MACHINES

Here the idea is to maintain the minimization of the flow time, but considering that the $\sum_j I\hat{Q}_j$ must be balanced across the available machines

- [1] Bertsimas, D., & Kallus, N. (2020). From predictive to prescriptive analytics. *Management Science*, 66(3), 1025-1044.
- [2] Mandi, J., Stuckey, P. J., & Guns, T. (2020, April). Smart predict-and-optimize for hard combinatorial optimization problems. In *Proceedings of the AAAI Conference on Artificial Intelligence* (Vol. 34, No. 02, pp. 1603-1610).
- [3] Li, Z., & Ierapetritou, M. (2008). Process scheduling under uncertainty: Review and challenges. *Computers & Chemical Engineering*, 32(4-5), 715-727.
- [4] Daniels, R. L., & Carrillo, J. E. (1997). β -Robust scheduling for single-machine systems with uncertain processing times. *IIE transactions*, 29(11), 977-985.
- [5] Magnusson, A., Punt, A. E., & Hilborn, R. (2013). Measuring uncertainty in fisheries stock assessment: the delta method, bootstrap, and MCMC. *Fish and Fisheries*, 14(3), 325-342.
- [6] Shrestha, D. L., & Solomatine, D. P. (2006). Machine learning approaches for estimation of prediction interval for the model output. *Neural Networks*, 19(2), 225-235.
- [7] Meinshausen, N., & Ridgeway, G. (2006). Quantile regression forests. *Journal of Machine Learning Research*, 7(6).
- [8] Rodríguez-Pérez, R., & Bajorath, J. (2019). Interpretation of compound activity predictions from complex machine learning models using local approximations and shapley values. *Journal of Medicinal Chemistry*, 63(16), 8761-8777.
- [10] Skutella, M., Sviridenko, M., & Uetz, M. (2016). Unrelated machine scheduling with stochastic processing times. *Mathematics of operations research*, 41(3), 851-864.



**INNOVACIÓN
INGENIERÍA UAI**
UNIVERSIDAD ADOLFO IBÁÑEZ

FACULTAD DE
INGENIERÍA Y CIENCIAS



SMART +
SUSTAINABLE

Using Quantile Forest For Robust Scheduling Of Astronomic Images Processing:

***Inform's Annual Meeting 2022
Indianapolis***

Gianfranco Speroni

gsperoni@alumnos.uai.cl

Luis Aburto

luis.aburto@uai.cl

Rodrigo Carrasco

rodrigo.carrascos@uai.cl

October 2022

	Models that estimate the processing time	Models that estimate a quantile of the processing time	Models that utilize the predictions of the quantile models to estimate the IQ
Descripción	They are models that calculate the conditioned expectation of the processing time given a set of observed features	They are models that calculate the conditioned quantile of the processing time given a set of observed features. In order to define an interval frequently 2 of these models are necessary	This model is a regular regression model of conditioned mean that is trained on the difference of the superior quantile prediction and the inferior quantile prediction. In other words on the IQ
Fórmula	$y_{pred} = E(Y = y X_1 = x_1, X_2 = x_2, \dots, X_n = x_n)$	$Q_{(Y X)}(\tau) = \inf\{y: F_{(Y X)}(y) \geq \tau\}$	$y_{pred} = E(Y = y X_1 = x_1, X_2 = x_2, \dots, X_n = x_n)$
Métricas	MAE, RMSE, MSE, MAPE, WMAPE and R2	Quantile Loss, absolute range, percentual range and data capture	MAE, RMSE, MSE, MAPE, WMAPE and R2



- Quantitative inspection through descriptive statistics
- Visual inspección through boxplots, scatterplots, histogramas y heatmaps
- Variable selection and outlier cleaning

Transform metadata

Box-Cox, Yeo-Johnson, Log, Sqrt, Categorical encoder y ordinal encoder

Train regression models

- Conditioned mean
 - Elastic Net, Kneighbors, Random Forest, ...
- Conditioned quantile
 - Quantile Forest, Quantile Gradient Boost, Quantile Light Gradient Boost y Bootstrap Methods
- Conditioned IQ
 - Gradient Boost, Ada Boost, Light Gradient Boost, ...

Select regression models

- Conditioned mean and IQ metrics
 - MAE, RMSE, MSE, MAPE, WMAPE y R2
- Conditioned quantile metrics
 - Quantile Loss, range, range [%] y data capture[%]

Formulate MIP for flow time minimization

- Oracle formulation
- Recommendation formulation
- Uncertainty penalization formulations

Evaluate the quality of the schedules

Mean gap[%] with respect to the oracle formulation

0 SELECTING A MODEL TO ESTIMATE THE IC(25%-75%)

Modelo	Tratamiento	Quantile loss (validation)	Quantile loss (training)	Percentil
Light Gradient Boosting	Standard scaler + log en target + log en features	6208	2317	25%
Ada Boost Bootstrap	Standard scaler + log en target	5552	NA	25%
Quantile Forest	NA	5671	2790	25%
Light Gradient Boosting	Standard scaler + log en features	8913	613	75%
Ada Boost Bootstrap	Standard scaler + log en target	7732	NA	75%
Quantile Forest	NA	8518	1366	75%