

# Compressible Tasks in Green Data Centers as Grid-Forming Support Assets

Based on work by

Anna Vandí, Ramon Aparicio-Pardo, Guillaume Urvoy-Keller.

*I3S, Université Côte d'Azur/CNRS*

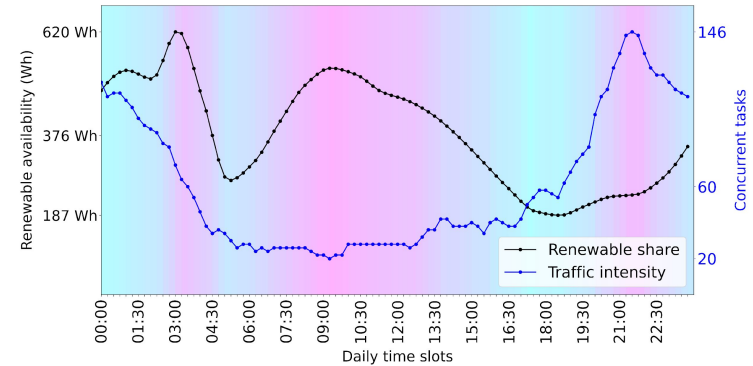
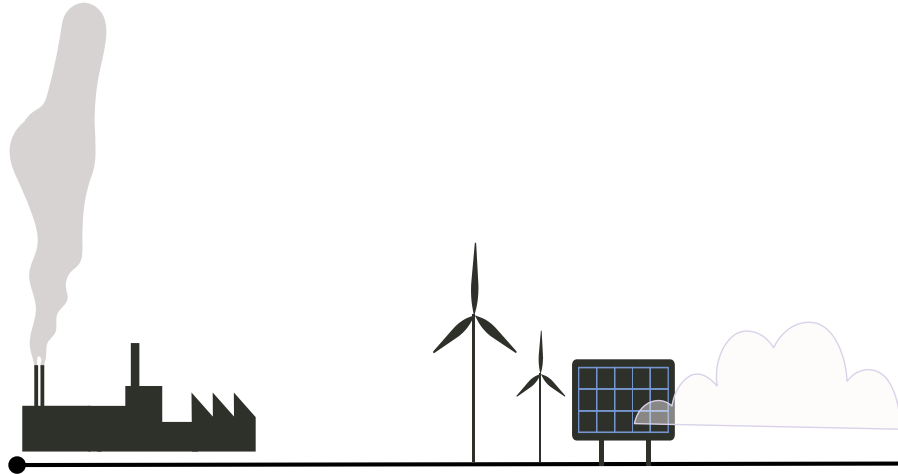
# Green Data-Centers equipped with microgrids

**Demand** in data centers **increases**

**Green Data-Centers**  
Equipped with  
**Microgrids**

**Uncertainty/Mismatch** challenge

- Renewable Power Generation [1]
- Load demand



# Renewable energy integration challenges grid stability

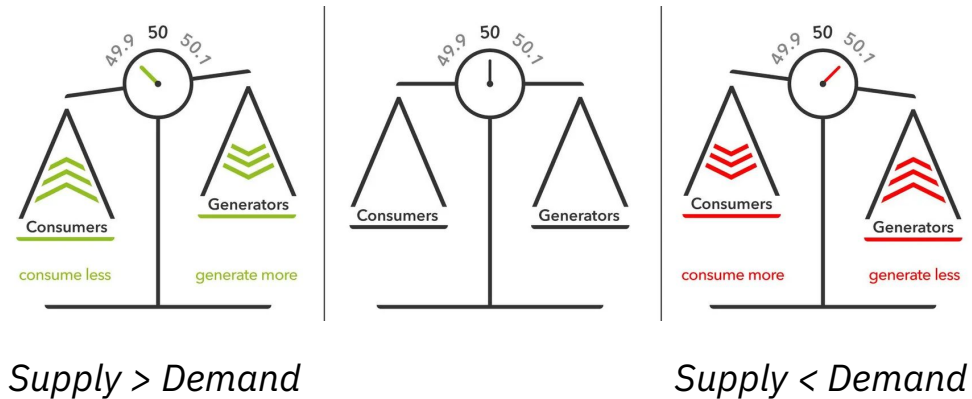
## Intermittency of Generation :

*supply/demand mismatch*

## Reduced Inertia:

*Conventional turbines provide natural inertia, stabilizing frequency<sup>[2]</sup>.*

### • Frequency (50 Hz) Deviations



[2] F. Ahmed et al., "Dynamic grid stability in low carbon power systems with minimum inertia," *Renewable Energy*, 2023. <https://www.sciencedirect.com/science/article/pii/S0960148123003774>

[3] ICS Investigation Expert Panel, "Grid Incident in Spain and Portugal on 28 April 2025: Factual Report," European Network of Transmission System Operators for Electricity (ENTSO-E), Tech. Rep., Oct. 3, 2025. source img: <https://www.next-kraftwerke.com/knowledge/afrr>

# Strategies to stabilize Renewable-Rich Grids

- Classical **hardware-based** solutions<sup>[2]</sup>
  - Energy Storage Systems (ESS)
  - Advanced Inverter Technologies as Grid-Forming Inverters (GFMI)
  
- Novel paradigm: **Demand-Side Flexibility**
  - Flexible Workloads as **Grid-Forming** Assets

[2] F. Ahmed et al., “Dynamic grid stability in low carbon power systems with minimum inertia,” Renewable Energy, 2023. <https://www.sciencedirect.com/science/article/pii/S0960148123003774>

[3] ICS Investigation Expert Panel, “Grid Incident in Spain and Portugal on 28 April 2025: Factual Report,” European Network of Transmission System Operators for Electricity (ENTSO-E), Tech. Rep., Oct. 3, 2025.

# Strategies to stabilize Renewable-Rich Grids

- Novel paradigm: ***Demand-Side Flexibility***
  - Flexible Workloads as **Grid-Forming** Assets

**Shift in space:** tasks across different servers/data centers

Drawbacks: latency, resource constraints, overhead <sup>[4],[6]</sup>

**Shift in time:** delaying or pausing/resuming workloads

Drawbacks: not always possible, more servers always ON <sup>[5]</sup>

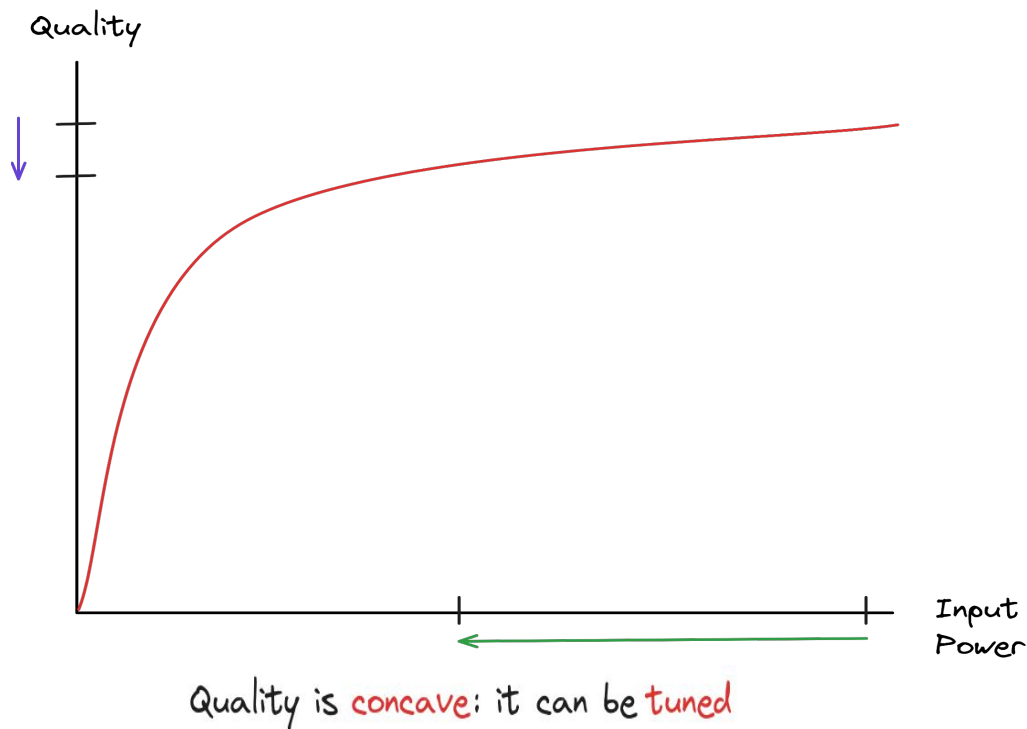
**PROPOSAL: “Shift” in size  
“compressible” tasks**

[4] W. E. Gribga et al., “Latency, energy and carbon aware collaborative resource allocation with consolidation and qos degradation strategies in edge computing,” in 2023 IEEE 29th ICPADS.

[5] B. Acun et al., “Carbon explorer: A holistic framework for designing carbon aware datacenters,” 2023.

[6] W. E. Gribga et al., “Energy-related impact of redefining self-consumption for distributed edge datacenters,” in 2024 IEEE 15th International Green and Sustainable Computing Conference (IGSC).

# “Compressible” tasks



## Video Transcoding Tasks

Transcoding output quality can be tuned [7]

## Image Classification Tasks

Inference optimization has impact on accuracy [8]

## LLM - Text Generation Tasks

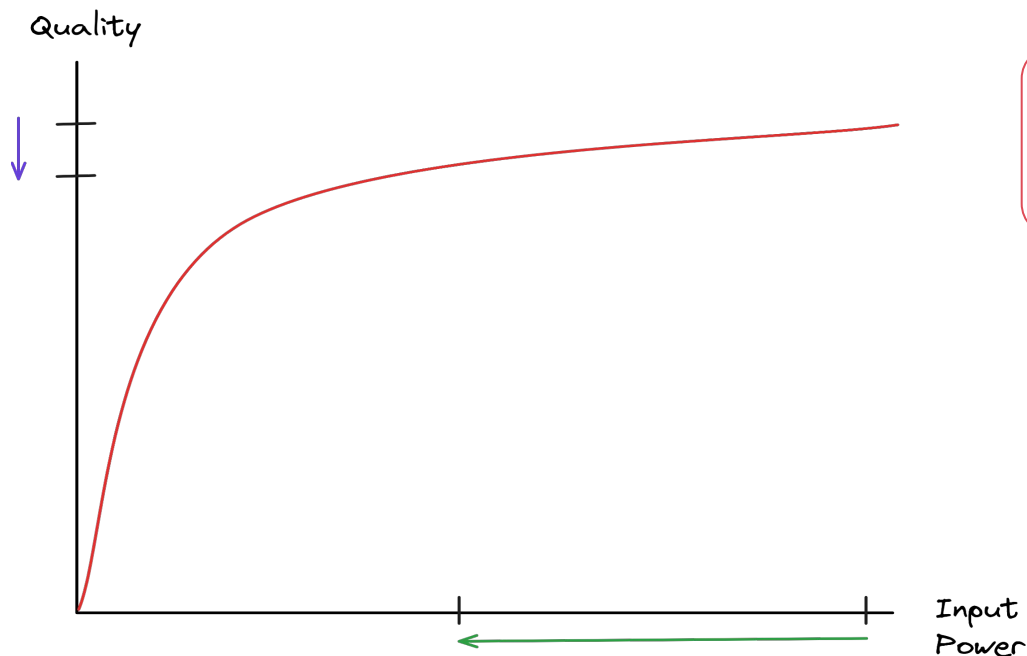
Inference optimization has impact on accuracy [9]

[7] Ramon Aparicio-Pardo, Karine Pires, Alberto Blanc, and Gwendal Simon. 2015. Transcoding live adaptive video streams at a massive scale in the cloud. (MMSys '15)

[8] H. Cai, C. Gan et al., “Once-for-all: Train one network and specialize it for efficient deployment,” in International Conference on Learning Representations, 2020.

[9] Xia, Heming, et al. "Tokenskip: Controllable chain-of-thought compression in llms." Proceedings of the 2025 Conference on Empirical Methods in Natural Language Processing. 2025.

# “Compressible” tasks



## Video Transcoding Tasks

Transcoding output quality can be tuned [7]

## Image Classification Tasks

Inference optimization has impact on accuracy [8]

## LLM - Text Generation Tasks

Inference optimization has impact on accuracy [9]

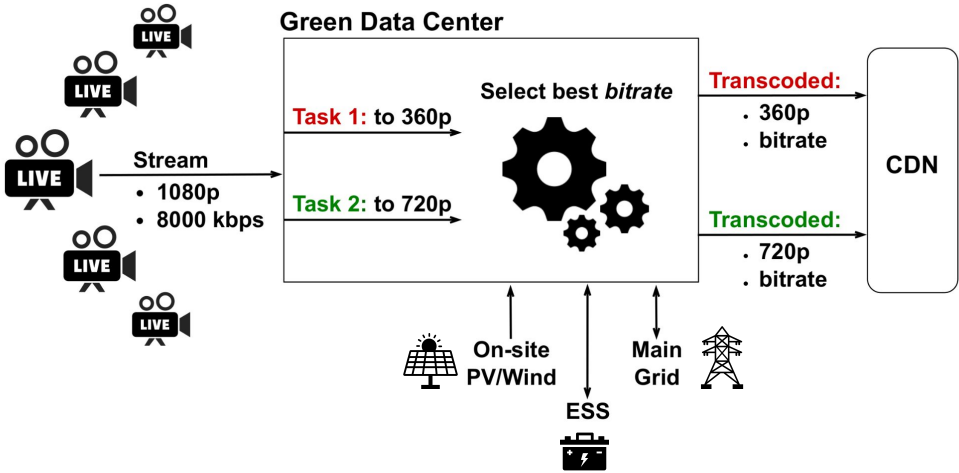
[7] Ramon Aparicio-Pardo, Karine Pires, Alberto Blanc, and Gwendal Simon. 2015. Transcoding live adaptive video streams at a massive scale in the cloud. (MMSys '15)

[8] H. Cai, C. Gan et al., “Once-for-all: Train one network and specialize it for efficient deployment,” in International Conference on Learning Representations, 2020.

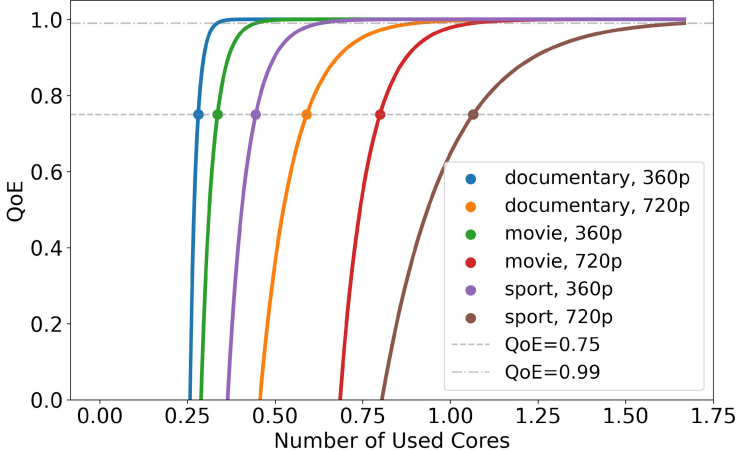
[9] Xia, Heming, et al. "Tokenskip: Controllable chain-of-thought compression in llms." Proceedings of the 2025 Conference on Empirical Methods in Natural Language Processing. 2025.

# A case study: live video transcoding

The System.



QoE for content type and resolution.



# Renewable Energy Aware Task Allocation (REATA)

## Oracle: Convex Formulation

$$\max_{\{\mathbf{x}\}} \sum_{t \in \mathcal{T}} \sum_{j \in \mathcal{J}} U_j(f_{tj} = F \cdot x_{tj})$$

s.t.

$$\sum_{j \in \mathcal{J}_i} P \cdot x_{tj} \leq \max(R_t, P_{ti}^{\min}),$$

$$x_{tj} \geq X_j^{\min},$$

$$x_{tj} \in \mathbb{R}_{\geq 0},$$

$$S_{ti} = \max(0, P_{ti}^{\min} - R_t)$$

---

$x_{tj} \in \mathbb{R}_{\geq 0}$	CPU usage by task $j$ at slot $t$ .
$U_j(x_j) \in \mathbb{R}_{\geq 0}$	Utility (QoE) of task $j$ at CPU usage $x_j$ .
$U^{\min} \in \mathbb{R}_{\geq 0}$	Minimal utility (QoE) level.
$X_j^{\min} \in \mathbb{R}_{\geq 0}$	Value for $x_j$ at minimal QoE level.
$P \in \mathbb{R}_{\geq 0}$	Maximum power consumption (in <i>Watts</i> ).
$R_t \in \mathbb{R}_{\geq 0}$	Renewable power ( <i>Watts</i> ) at slot $t$ .
$S_{ti} \in \mathbb{R}_{\geq 0}$	Power <i>imported</i> at slot-interval pair $(t, i)$ .
$F \in \mathbb{R}_{\geq 0}$	CPU core frequency (in <i>GHz</i> ).

---

## ONLINE-REATA Algorithm

---

**Algorithm 1:** ONLINE-REATA: abstract task allocation

---

**Data:**  $\mathcal{J}_{i-1}, \mathcal{H}_{i-1}, \mathcal{L}_{i-1}, \tilde{N}_i, S_{i-1}, R$

$k$ : arriving task at interval  $i$  (slot  $t$ ).

**Result:** Optimal allocation  $x_k$  for task  $k$

```

1 for each arriving task  $k$  do
2   if  $|\mathcal{J}_{i-1}| + 1 \geq \tilde{N}_i$  then
3      $\tilde{\mathcal{H}}_i \leftarrow \mathcal{H}_{i-1} \cup \{k\}$ ;
4   else
5     Generate artificial tasks  $\mathcal{A}_{it}$  based on  $\tilde{N}_i$ ;
6      $\tilde{\mathcal{H}}_i \leftarrow \mathcal{H}_{i-1} \cup \{k\} \cup \mathcal{A}_{it}$ ;
7   Compute  $\lambda$  using (8) with  $(\tilde{\mathcal{H}}_i, \mathcal{L}_{i-1}, R, S_{i-1})$ ;
8   if  $\lambda > U'_k(X_k^{\min})$  then
9      $x_k \leftarrow X_k^{\min}$ ;
10     $\mathcal{L}_i \leftarrow \mathcal{L}_{i-1} \cup \{k\}$ ;
11  else
12     $x_k \leftarrow -\frac{1}{\theta_k F} \ln\left(\frac{\lambda}{\tau_k \theta_k P}\right)$ ;
13     $\mathcal{H}_i \leftarrow \mathcal{H}_{i-1} \cup \{k\}$ ;
14   $\mathcal{J}_i \leftarrow \mathcal{H}_i \cup \mathcal{L}_i$ ;

```

---

# Renewable Energy Aware Task Allocation (REATA)

## Oracle: Convex Formulation

$$\max_{\{x\}} \sum_{t \in \mathcal{T}} \sum_{j \in \mathcal{J}} U_j(f_{tj} = F \cdot x_{tj}) \quad \text{Max Sum Utilities}$$

s.t.

$$\sum_{j \in \mathcal{J}_i} P \cdot x_{tj} \leq \max(R_t, P_{ti}^{\min}), \quad \text{Power Budget}$$

$$x_{tj} \geq X_j^{\min},$$

Minimum Allocation  
(QoE)

$$x_{tj} \in \mathbb{R}_{\geq 0},$$

$$S_{ti} = \max(0, P_{ti}^{\min} - R_t)$$

$x_{tj}$  decision variable

$x_{tj} \in \mathbb{R}_{>0}$	CPU usage by task $j$ at slot $t$ .
$U_j(x_j) \in \mathbb{R}_{\geq 0}$	Utility (QoE) of task $j$ at CPU usage $x_j$ .
$U^{\min} \in \mathbb{R}_{\geq 0}$	Minimal utility (QoE) level.
$X_j^{\min} \in \mathbb{R}_{\geq 0}$	Value for $x_j$ at minimal QoE level.
$P \in \mathbb{R}_{>0}$	Maximum power consumption (in <i>Watts</i> )
$R_t \in \mathbb{R}_{>0}$	Renewable power ( <i>Watts</i> ) at slot $t$ .
$S_{ti} \in \mathbb{R}_{\geq 0}$	Power <i>imported</i> at slot-interval pair $(t, i)$ .
$F \in \mathbb{R}_{>0}$	CPU core frequency (in <i>GHz</i> ).

## ONLINE-REATA Algorithm

Algorithm 1: ONLINE-REATA: abstract task allocation

**Data:**  $\mathcal{J}_{i-1}, \mathcal{H}_{i-1}, \mathcal{L}_{i-1}, \tilde{N}_i, S_{i-1}, R$   
 **$k$ :** arriving task at interval  $i$  (slot  $t$ ).  
**Result:** Optimal allocation  $x_k$  for task  $k$

```

1 for each arriving task  $k$  do
2   if  $|\mathcal{J}_{i-1}| + 1 \geq \tilde{N}_i$  then
3      $\tilde{\mathcal{H}}_i \leftarrow \mathcal{H}_{i-1} \cup \{k\}$ ;
4   else
5     Generate artificial tasks  $\mathcal{A}_{it}$  based on  $\tilde{N}_i$ ;
6      $\tilde{\mathcal{H}}_i \leftarrow \mathcal{H}_{i-1} \cup \{k\} \cup \mathcal{A}_{it}$ ;
7   Compute  $\lambda$  using (8) with  $(\tilde{\mathcal{H}}_i, \mathcal{L}_{i-1}, R, S_{i-1})$ ;
8   if  $\lambda > U'_k(X_k^{\min})$  then
9      $x_k \leftarrow X_k^{\min}$ ;
10     $\mathcal{L}_i \leftarrow \mathcal{L}_{i-1} \cup \{k\}$ ;
11  else
12     $x_k \leftarrow -\frac{1}{\theta_k F} \ln\left(\frac{\lambda}{\tau_k \theta_k P}\right)$ ;
13     $\mathcal{H}_i \leftarrow \mathcal{H}_{i-1} \cup \{k\}$ ;
14   $\mathcal{J}_i \leftarrow \mathcal{H}_i \cup \mathcal{L}_i$ ;

```

# Renewable Energy Aware Task Allocation (REATA)

## Oracle: Convex Formulation

$$\begin{aligned} & \max_{\{\mathbf{x}\}} \sum_{t \in \mathcal{T}} \sum_{j \in \mathcal{J}} U_j(f_{tj} = F \cdot x_{tj}) \\ & \text{s.t.} \\ & \sum_{j \in \mathcal{J}_i} P \cdot x_{tj} \leq \max(R_t, P_{ti}^{\min}), \\ & x_{tj} \geq X_j^{\min}, \\ & x_{tj} \in \mathbb{R}_{\geq 0}, \\ & S_{ti} = \max(0, P_{ti}^{\min} - R_t) \end{aligned}$$

---

$x_{tj} \in \mathbb{R}_{\geq 0}$	CPU usage by task $j$ at slot $t$ .
$U_j(x_j) \in \mathbb{R}_{\geq 0}$	Utility (QoE) of task $j$ at CPU usage $x_j$ .
$U^{\min} \in \mathbb{R}_{\geq 0}$	Minimal utility (QoE) level.
$X_j^{\min} \in \mathbb{R}_{\geq 0}$	Value for $x_j$ at minimal QoE level.
$P \in \mathbb{R}_{\geq 0}$	Maximum power consumption (in <i>Watts</i> ).
$R_t \in \mathbb{R}_{\geq 0}$	Renewable power ( <i>Watts</i> ) at slot $t$ .
$S_{ti} \in \mathbb{R}_{\geq 0}$	Power <i>imported</i> at slot-interval pair $(t, i)$ .
$F \in \mathbb{R}_{\geq 0}$	CPU core frequency (in <i>GHz</i> ).

---

## ONLINE-REATA Algorithm

---

**Algorithm 1:** ONLINE-REATA: abstract task allocation

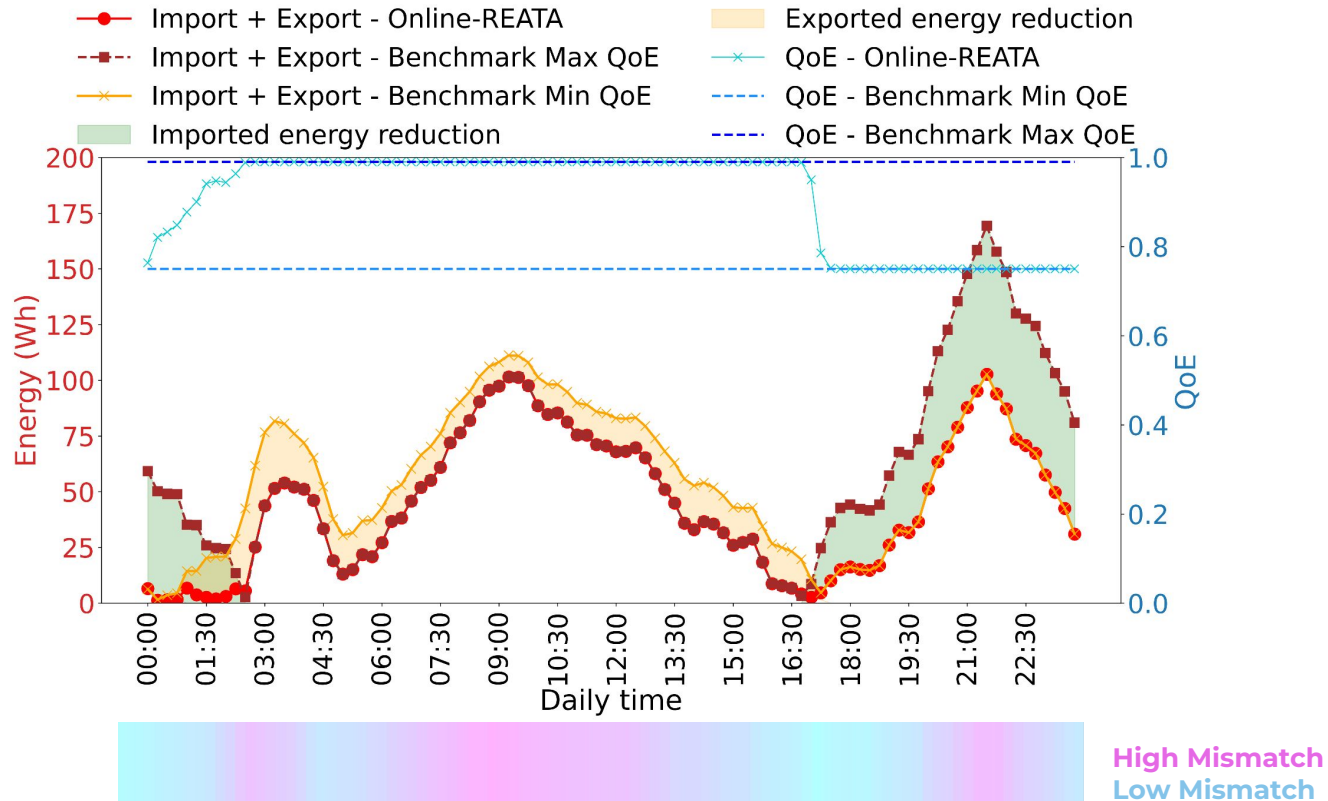
---

**Data:**  $\mathcal{J}_{i-1}, \mathcal{H}_{i-1}, \mathcal{L}_{i-1}, \hat{N}_i, S_{i-1}, R$   
 **$k$ :** arriving task at interval  $i$  (slot  $t$ ).  
**Result:** Optimal allocation  $x_k$  for task  $k$

- 1 **for each arriving task  $k$  do**
- 2     **if**  $|\mathcal{J}_{i-1}| + 1 \geq \hat{N}_i$  **then**
- 3          $\hat{\mathcal{H}}_i \leftarrow \mathcal{H}_{i-1} \cup \{k\}$ ;
- 4     **else**
- 5         Generate artificial tasks  $\mathcal{A}_{iI}$  based on  $\hat{N}_i$ ;
- 6          $\hat{\mathcal{H}}_i \leftarrow \mathcal{H}_{i-1} \cup \{k\} \cup \mathcal{A}_{iI}$ ;
- 7         Compute  $\lambda$  using (8) with  $(\hat{\mathcal{H}}_i, \mathcal{L}_{i-1}, R, S_{i-1})$ ;
- 8         **if**  $\lambda > U'_k(X_k^{\min})$  **then**
- 9              $x_k \leftarrow X_k^{\min}$ ;
- 10             $\mathcal{L}_i \leftarrow \mathcal{L}_{i-1} \cup \{k\}$ ;
- 11         **else**
- 12              $x_k \leftarrow -\frac{1}{\theta_k F} \ln\left(\frac{\lambda}{r_k \theta_k F}\right)$ ;
- 13              $\hat{\mathcal{H}}_i \leftarrow \mathcal{H}_{i-1} \cup \{k\}$ ;
- 14          $\mathcal{J}_i \leftarrow \mathcal{H}_i \cup \mathcal{L}_i$ ;

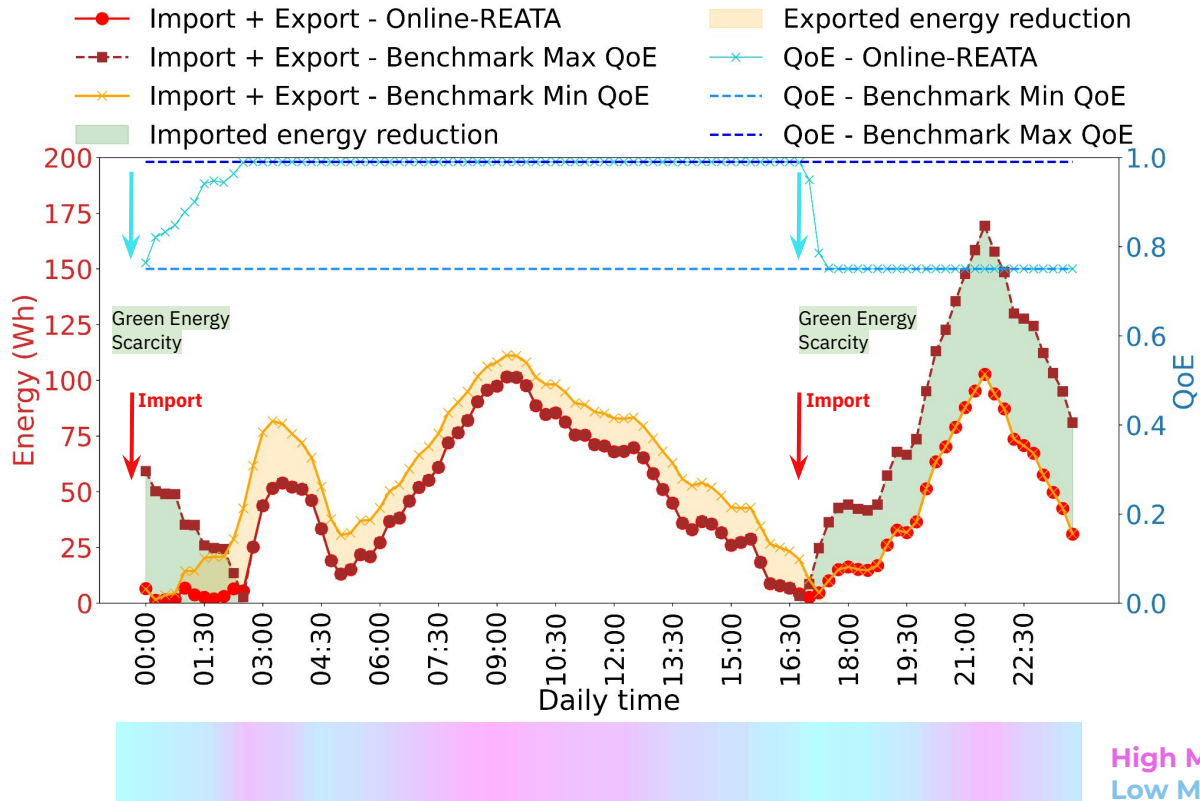
---

# Some results:

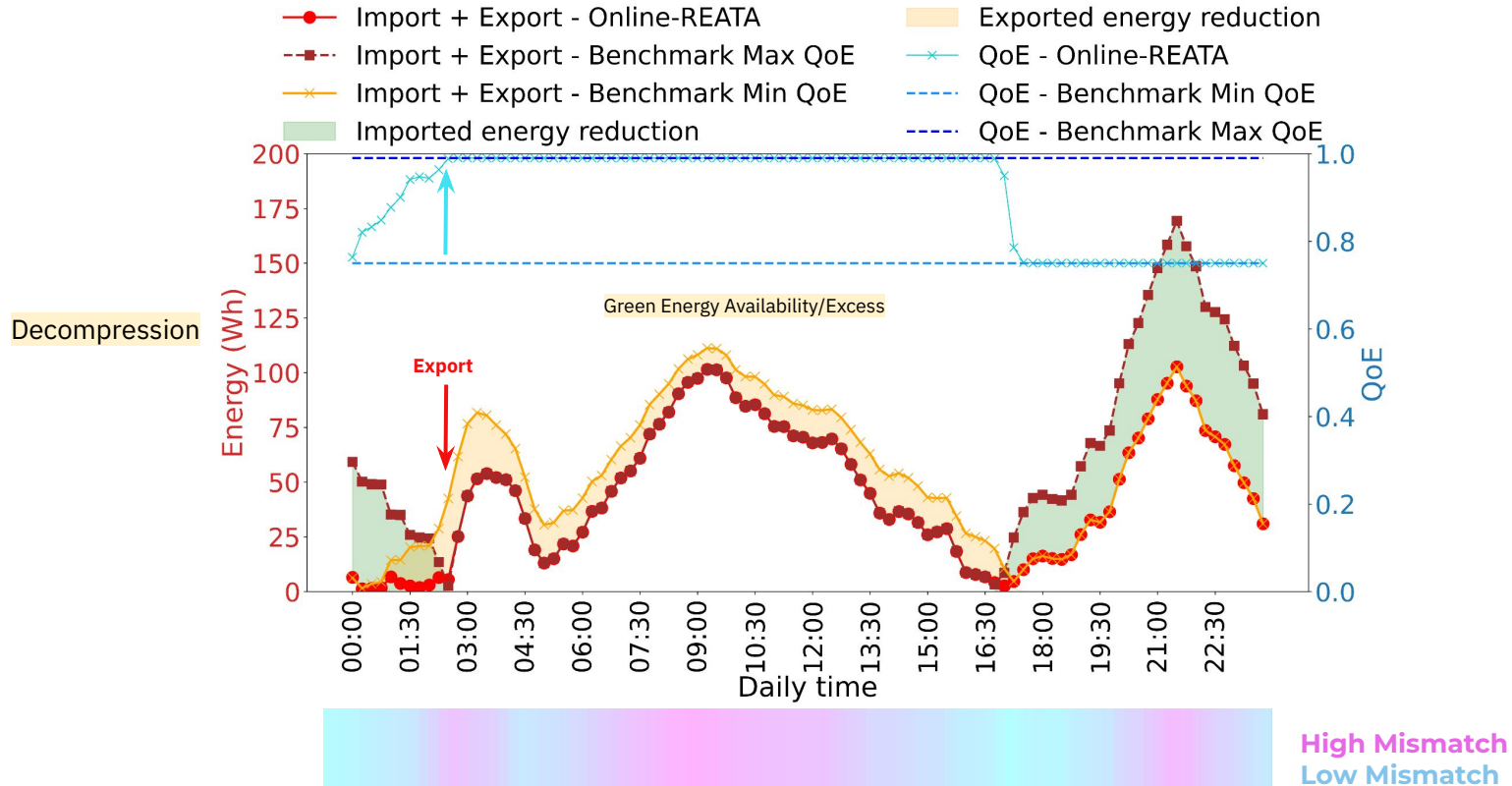


# Some results:

Compression



# Some results:



# Some results

## Self-Sufficiency Rate (SSR) [11]

- share of load covered by renewables

$$SSR_{PV+ESS} = \frac{PV2load + ESS_{PV}2load}{PV2load + ESS2load + Grid2load}$$

- 📌 High SSR  $\Rightarrow$  less import from the main grid

## Self-Consumption Rate (SCR) [11]

- share of renewable energy used over total available

$$SCR_{PV+ESS} = \frac{PV2load + ESS_{PV}2load}{PV2grid + PV2load + PV2ESS}$$

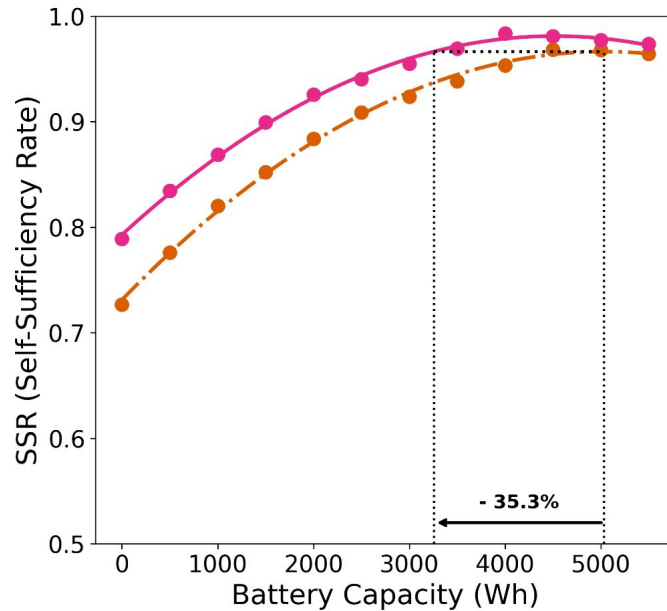
- 📌 High SCR  $\Rightarrow$  less export to the main grid

# Some results

## Self-Sufficiency Rate (SSR) [11]

- share of load covered by renewables

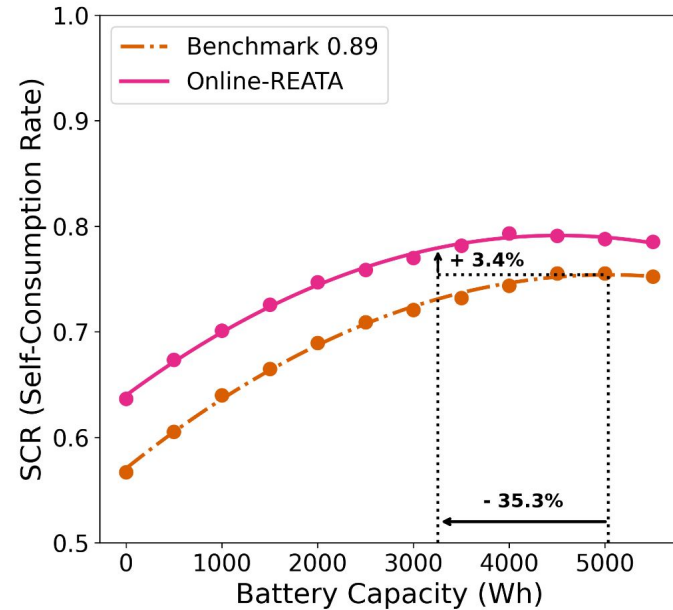
$$SSR_{PV+ESS} = \frac{PV2load + ESS_{PV}2load}{PV2load + ESS2load + Grid2load}$$



## Self-Consumption Rate (SCR) [11]

- share of renewable energy used over total available

$$SCR_{PV+ESS} = \frac{PV2load + ESS_{PV}2load}{PV2grid + PV2load + PV2ESS}$$

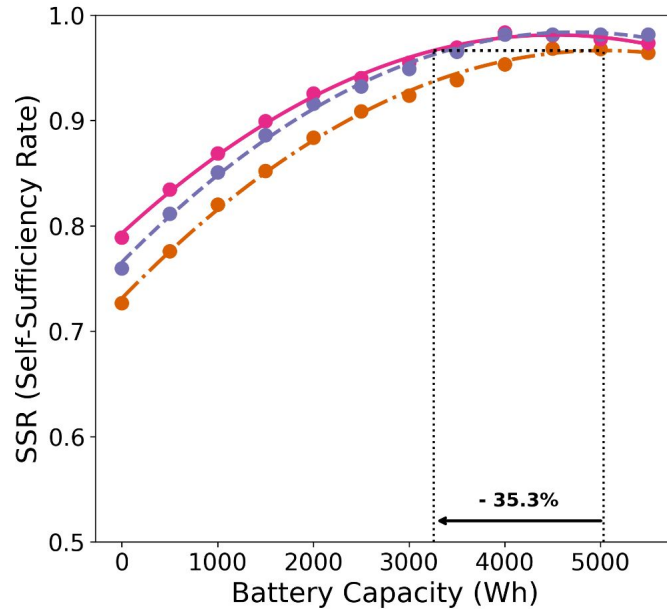


# Some results

## Self-Sufficiency Rate (SSR) [11]

- share of load covered by renewables

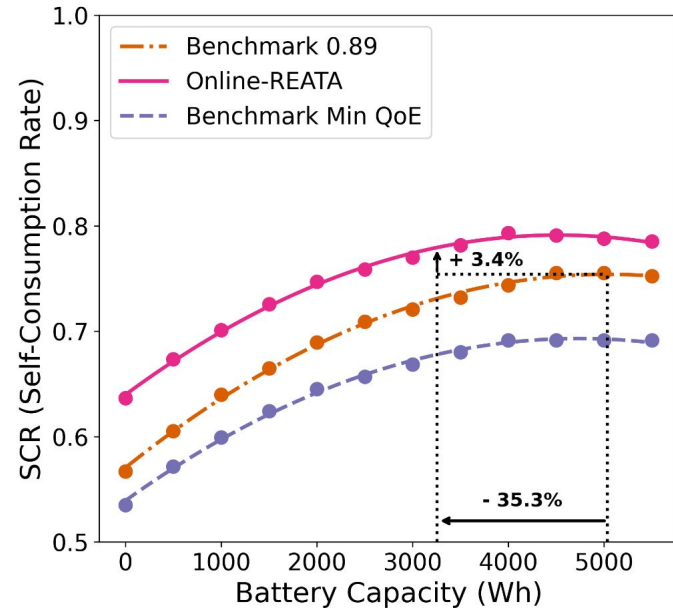
$$SSR_{PV+ESS} = \frac{PV2load + ESS_{PV}2load}{PV2load + ESS2load + Grid2load}$$



## Self-Consumption Rate (SCR) [11]

- share of renewable energy used over total available

$$SCR_{PV+ESS} = \frac{PV2load + ESS_{PV}2load}{PV2grid + PV2load + PV2ESS}$$

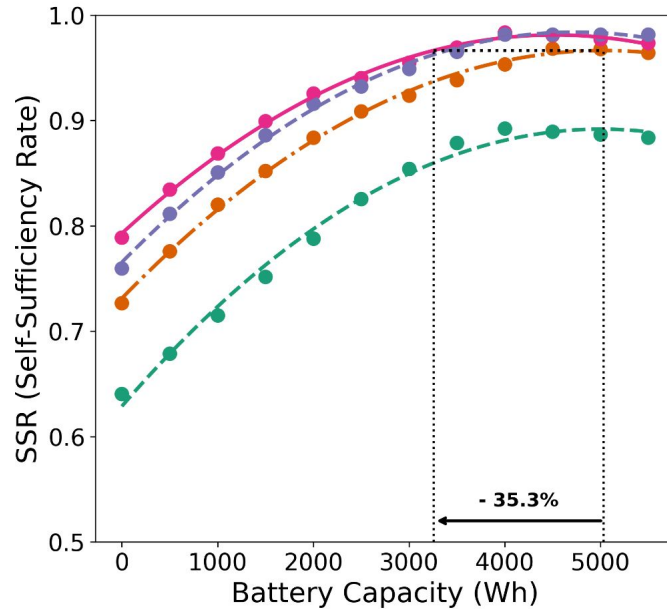


# Some results

## Self-Sufficiency Rate (SSR) [11]

- share of load covered by renewables

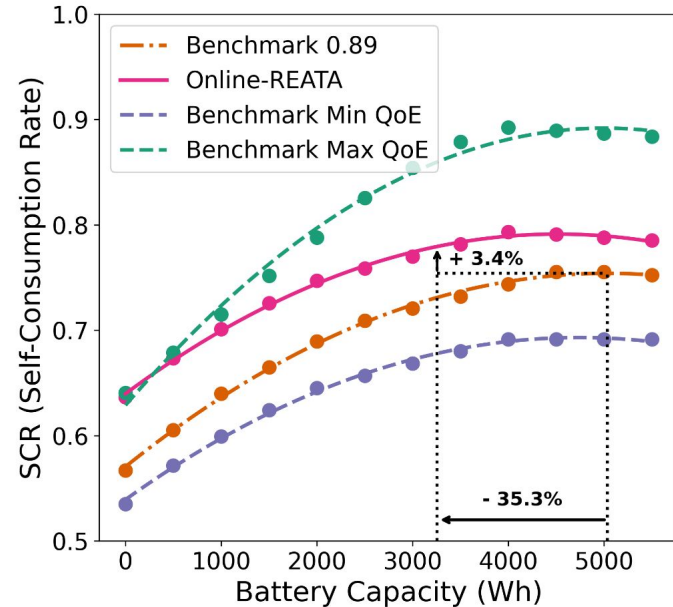
$$SSR_{PV+ESS} = \frac{PV2load + ESS_{PV}2load}{PV2load + ESS2load + Grid2load}$$



## Self-Consumption Rate (SCR) [11]

- share of renewable energy used over total available

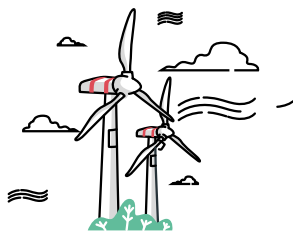
$$SCR_{PV+ESS} = \frac{PV2load + ESS_{PV}2load}{PV2grid + PV2load + PV2ESS}$$



# Future Directions:



- Extension to other **compressible tasks** (RLMs and CoT pruning)
- Model selection (Object detection, text generation, code generation,...)
- Combination with **time and space shifts** (batches as GFM asset, servers consolidation)
- **Distributed** environments integration
- Social Incentives



Thank You!



[vandi@i3s.unice.fr](mailto:vandi@i3s.unice.fr)