# Caching or pre-fetching? The role of hazard rates.

Andres Ferragut
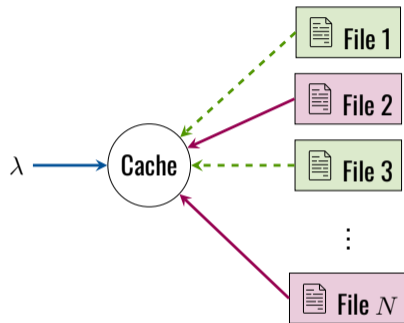
joint work with Matias Carrasco and Fernando Paganini

Universidad ORT Uruguay

# The caching problem



- Consider a local memory system that handles items from a catalog of $N$ objects.

- Requests for objects arrive as a random process.

- The memory (cache) can locally store $C < N$ of them.

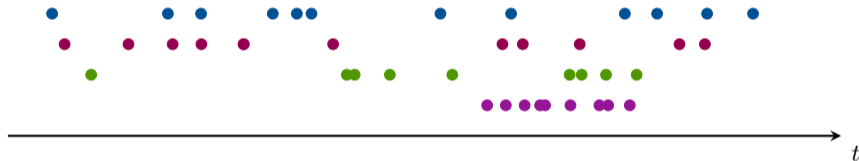- If item is in cache, we have a hit. Otherwise, it is a miss.

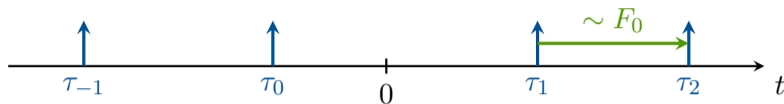Objective: for a given arrival stream, maximize the steady-state hit rate.

Introduced in [Fofack et al. 2014]

- Assume requests for item $i$ come from a point process of intensity $\lambda_i$ (popularities).



- At each point in time we must decide which items must be stored locally.

# Two important distributions:



- Inter-arrival distribution: Typical distance between points...

$$\tau_{k+1} - \tau_k \sim F_0(t), \quad E[\tau_{k+1} - \tau_k] = 1/\lambda$$

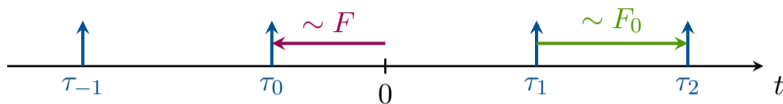# Two important distributions:



- Inter-arrival distribution: Typical distance between points...

$$\tau_{k+1} - \tau_k \sim F_0(t), \quad E[\tau_{k+1} - \tau_k] = 1/\lambda$$

- Age distribution: Distance from the last point in the current interval (sampling bias)!

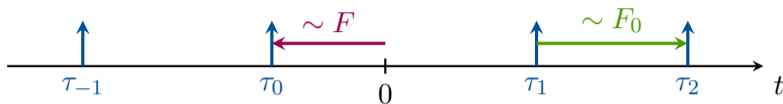$$F(t) := \lambda \int_0^t 1 - F_0(s)ds,$$

# Two important distributions:



■ **Inter-arrival distribution:** Typical distance between points...

$$\tau_{k+1} - \tau_k \sim F_0(t), \quad E[\tau_{k+1} - \tau_k] = 1/\lambda$$

■ **Age distribution:** Distance from the last point in the current interval (sampling bias)!

$$F(t) := \lambda \int_0^t 1 - F_0(s)ds,$$

Note: you can formalize this under the Palm probability framework for stationary point processes.

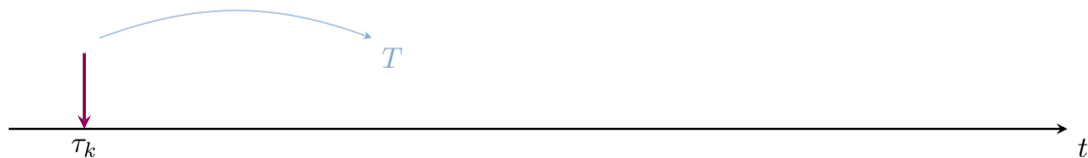- Upon request arrival for item $i$, check for presence.

$\tau_k$

$t$

# Populating a cache: timer based (TTL) policies

- Upon request arrival for item $i$, check for presence.
- If new, store item and start a timer $T_i$ to evict.



$T$

$\tau_k$

$t$

# Populating a cache: timer based (TTL) policies

- Upon request arrival for item $i$, check for presence.
- If new, store item and start a timer $T_i$ to evict.
- If present, reset timer to $T_i$.

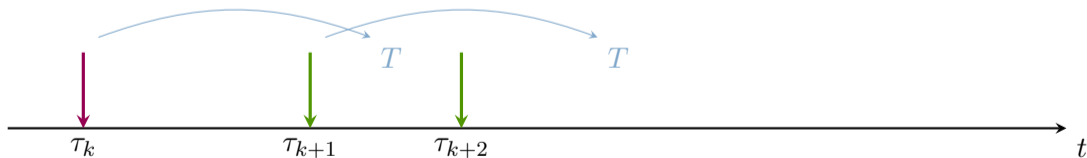# Populating a cache: timer based (TTL) policies

- Upon request arrival for item $i$, check for presence.
- If new, store item and start a timer $T_i$ to evict.
- If present, reset timer to $T_i$.
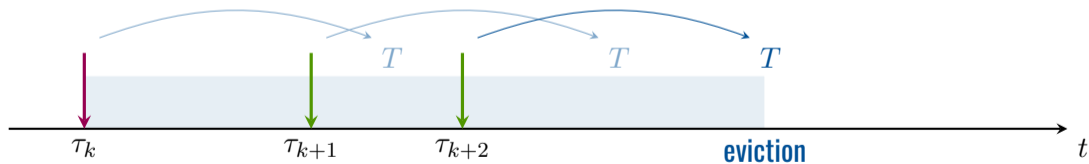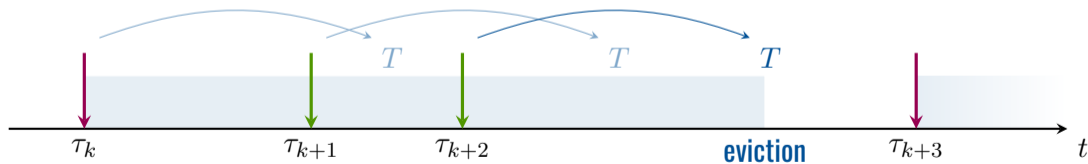- Upon timer expiration, evict the content.

# Populating a cache: timer based (TTL) policies

- Upon request arrival for item $i$, check for presence.
- If new, store item and start a timer $T_i$ to evict.
- If present, reset timer to $T_i$.
- Upon timer expiration, evict the content.
- Keep timers $T_i$ such that average cache occupation is $C$.

# Structure of the optimal caching policy

- The crucial magnitude is the hazard rate of $F_0$:

$$\eta(t) := \frac{f_0(t)}{1 - F_0(t)}$$

- Likelihood of a request at time $t$, given the current interval has age $t$.

# Structure of the optimal caching policy

- The crucial magnitude is the hazard rate of $F_0$:

$$\eta(t) := \frac{f_0(t)}{1 - F_0(t)}$$

- Likelihood of a request at time $t$, given the current interval has age $t$.

## Theorem (F', Rodriguez, Paganini – 2016)

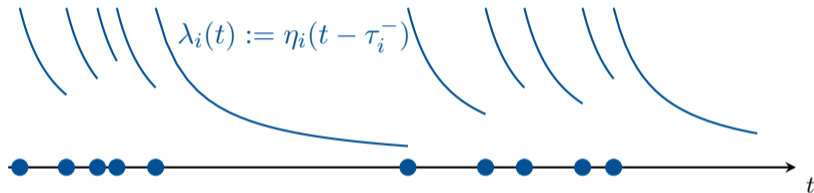If the $F_0^{(i)}$ have decreasing hazard rates, then the optimal TTL policy satisfies:

$$\eta_i(T_i^*) \geqslant \theta^*,$$

whenever $T_i^* > 0$ (i.e. the item is cached).
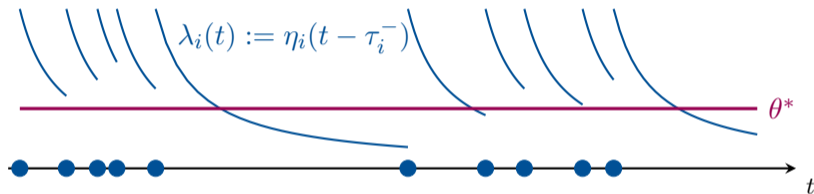Moreover, inequality is strict iff $T_i^* = \infty$ (item always stored).

Decreasing hazard rates corresponds to bursty traffic:



$$\lambda_i(t) := \eta_i(t - \tau_i^-)$$

# Why caching helps in this case?

Decreasing hazard rates corresponds to bursty traffic:



$$\lambda_i(t) := \eta_i(t - \tau_i^-)$$

$\theta^*$

- An arrival makes a subsequent arrival more likely.
- Store it while its likelihood is high enough (above a threshold).

# Why caching helps in this case?

Decreasing hazard rates corresponds to bursty traffic:



$$\lambda_i(t) := \eta_i(t - \tau_i^-)$$

$\theta^*$

- An arrival makes a subsequent arrival more likely.
- Store it while its likelihood is high enough (above a threshold).

# What about other types of traffic?
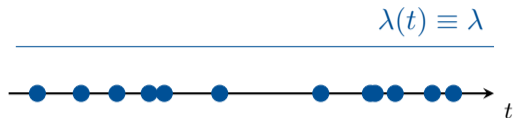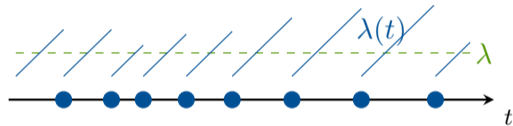


Constant hazard rate $\rightarrow$ Poisson process.

Increasing hazard rate $\rightarrow$ more periodic!

# What about other types of traffic?



$\lambda(t) \equiv \lambda$

Constant hazard rate $\rightarrow$ Poisson process.

$\lambda(t)$ $\lambda$

Increasing hazard rate $\rightarrow$ more periodic!

Theorem: for these types of traffic, keep the most popular is the optimal caching policy.

$\lambda(t) \equiv \lambda$

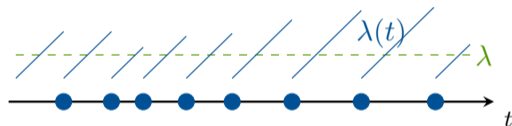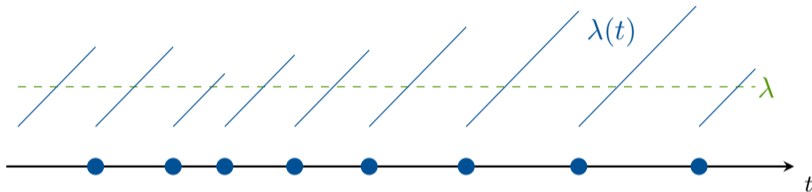Constant hazard rate $\rightarrow$ Poisson process.

$\lambda(t)$ $\lambda$

Increasing hazard rate $\rightarrow$ more periodic!

Theorem: for these types of traffic, keep the most popular is the optimal caching policy.

Can we improve upon this?

- Once you have seen a request, it's less likely to see the same item again for a while.



- What is the timer based equivalent of this case?

# Thinking about increasing hazard rates...
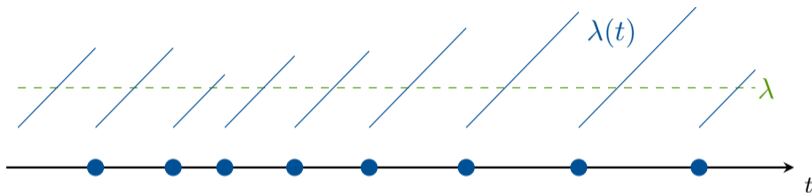
- Once you have seen a request, it's less likely to see the same item again for a while.



- What is the timer based equivalent of this case?

## Key insight

The question now is not how long we should remember something, but instead how long we should forget about it!

# Timer based pre-fetching policy

- Upon request arrival for item $i$, check for presence.

$\tau_k$

$t$

- Upon request arrival for item $i$, check for presence.
- If not-present: start a timer $T_i$.

# Timer based pre-fetching policy

- Upon request arrival for item $i$, check for presence.
- If not-present: start a timer $T_i$.
- If present: remove content and reset timer to $T_i$.

# Timer based pre-fetching policy

- Upon request arrival for item $i$, check for presence.
- If not-present: start a timer $T_i$.
- If present: remove content and reset timer to $T_i$.
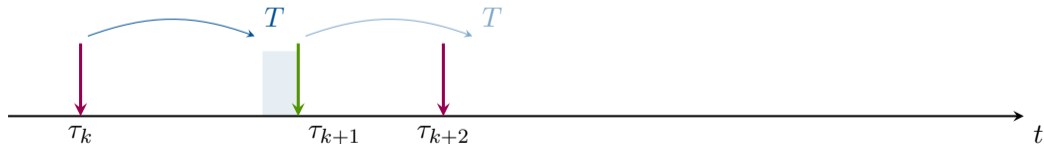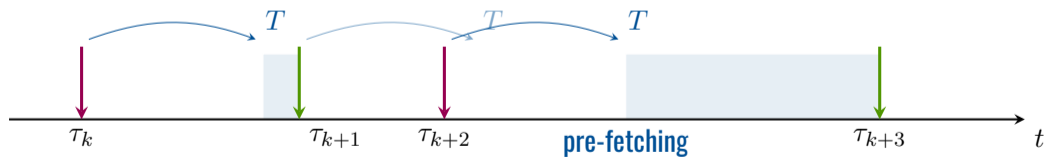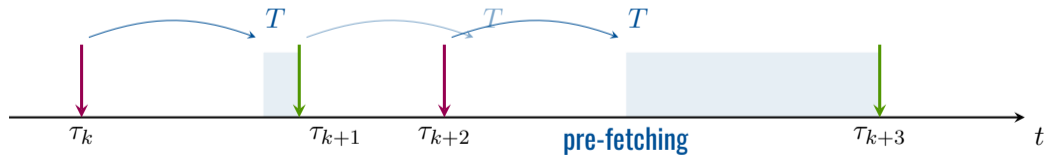- Upon timer expiration, pre-fetch the content.

# Timer based pre-fetching policy

- Upon request arrival for item $i$, check for presence.
- If not-present: start a timer $T_i$.
- If present: remove content and reset timer to $T_i$.
- Upon timer expiration, pre-fetch the content.
- Keep timers $T_i$ such that average cache occupation is $C$.

# Timer based pre-fetching

Consider a single item with a timer $T$ and its request process:

**Hit probability:** next arrival occurs after timer expires.

**Occupation probability:** probability that timer has expired by $0$ since last arrival.



Hit probability $= 1 - F_0(T)$

Avg. occupation $= 1 - F(T)$

# Choosing the optimal timers

Requests come from independent sources with intensities $\lambda_i$ and inter-arrival distribution $F_i$:

## Problem (Optimal pre-fetching policy)

Choose timers $T_i \geqslant 0$ such that:

$$\max_{T_i \geqslant 0} \sum_i \lambda_i (1 - F_0^{(i)}(T_i))$$

subject to:

$$\sum_i (1 - F^{(i)}(T_i)) \leqslant C$$

# Choosing the optimal timers

Requests come from independent sources with intensities $\lambda_i$ and inter-arrival distribution $F_i$:

## Problem (Optimal pre-fetching policy)

Choose timers $T_i \geqslant 0$ such that:

$$\min_{T_i \geqslant 0} \sum_i \lambda_i F_0^{(i)}(T_i)$$

subject to:

$$\sum_i F^{(i)}(T_i) \geqslant N - C$$

# Choosing the optimal timers
Change of variables

- Apply the change of variables $u_i = F^{(i)}(T_i)$.

- Note that $u_i$ is the probability of not being stored.

- The problem becomes:

$$\min_{u_i \in [0,1]} \sum_i \lambda_i \left[ F_0^{(i)} \circ (F^{(i)})^{-1} \right] (u_i)$$

subject to:

$$\sum_i u_i \geqslant N - C$$

# Choosing the optimal timers

Lagrangian duality

- Objective gradient:

$$\frac{\partial}{\partial u_i} \lambda_i F^{(i)} \circ (F^{(i)})^{-1}(u_i) = \frac{\lambda_i f_0^{(i)}((F^{(i)})^{-1}(u_i))}{\lambda_i \left(1 - F_0^{(i)}((F^{(i)})^{-1}(u_i))\right)} = \eta_i((F^{(i)})^{-1}(u_i))$$

- Increasing! $\rightarrow$ Proper convex optimization problem.

# Choosing the optimal timers

Lagrangian duality

- Objective gradient:

$$\frac{\partial}{\partial u_i} \lambda_i F^{(i)} \circ (F^{(i)})^{-1}(u_i) = \frac{\lambda_i f_0^{(i)}((F^{(i)})^{-1}(u_i))}{\lambda_i \left(1 - F_0^{(i)}((F^{(i)})^{-1}(u_i))\right)} = \eta_i((F^{(i)})^{-1}(u_i))$$

- Increasing! → Proper convex optimization problem.
- Lagrangian duality:

$$\mathcal{L}(u, \theta) = \sum_{i=1}^{N} \lambda_i F_0^{(i)} \left((F^{(i)})^{-1}(u_i)\right) + \theta \left(N - C - \sum_{i=1}^{N} u_i\right)$$

$$= \sum_{i=1}^{N} \left[\lambda_i F_0^{(i)} \left((F^{(i)})^{-1}(u_i)\right) - \theta u_i\right] + \theta(N - C).$$

# Optimal Timer-based pre-fetching

## Theorem

If the $F_0^{(i)}$ satisfy the IHR property, there exists a unique threshold $\theta^* \geqslant 0$ such that the optimal timers satisfy:

$$\eta_i(T_i^*) \geqslant \theta^*,$$

whenever $T_i^* < \infty$ (pre-fetching).
The inequality is strict if and only if $T_i^* = 0$, i.e. the content is always stored.

# Optimal Timer-based pre-fetching

## Theorem

If the $F_0^{(i)}$ satisfy the IHR property, there exists a unique threshold $\theta^* \geqslant 0$ such that the optimal timers satisfy:
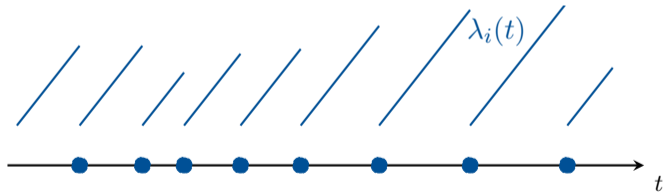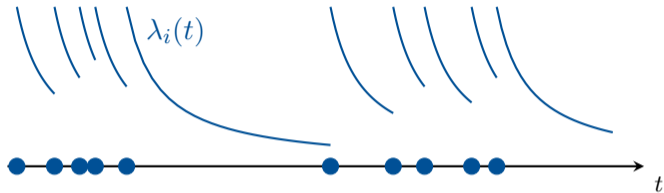
$$\eta_i(T_i^*) \geqslant \theta^*,$$

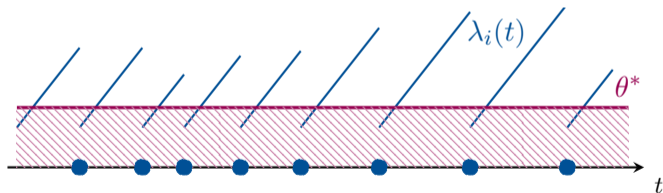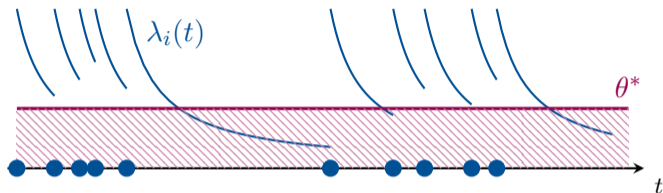whenever $T_i^* < \infty$ (pre-fetching).
The inequality is strict if and only if $T_i^* = 0$, i.e. the content is always stored.

Remark: The policy is also a threshold policy, like the caching case.
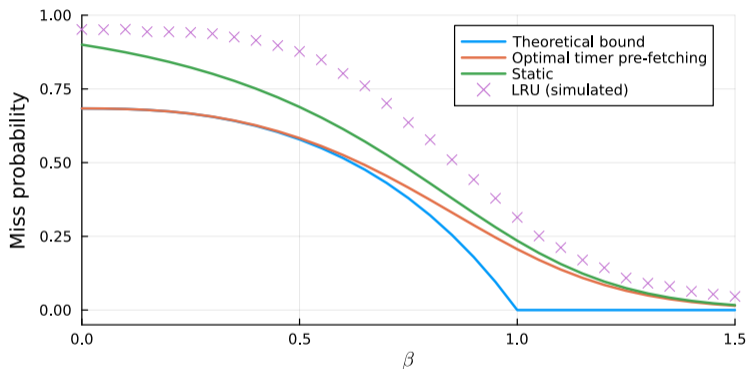
# A tale of two policies...



Both policies are just the same policy!

- Keep a hazard rate threshold $\theta$ for storing a content
- Compute $\theta^*$ such that avg. memory occupation is $C$.

# Simulation example

Erlang ($k = 5$) inter-arrival times, Zipf $\propto n^{-\beta}$ popularities, varying $\beta$, $N = 10000$, $C = 1000$.



- Optimal pre-fetching improves over the static policy.
- Classical caching (e.g. LRU) is a very bad idea for regular traffic.

## Theorem

Consider a family of local memory systems, indexed by $N$, with inter-request times coming from a common scale family, and memory size $C_N = cN$. Then the hazard rate threshold $\theta_N^*$ verifies:

$$\theta_N \xrightarrow[N \to \infty]{} \theta^*,$$

where $\theta^*$ is the solution of:

$$G_\infty(\theta^*) = 1 - c,$$

and $G_\infty$ depends on the age distribution $F$ and the popularity distribution $L$.

# Main results
Asymptotic miss probability

## Theorem

Consider a family of local memory systems as before, then the miss probability for system $N$, $M_N$, satisfies:

$$M_N \underset{N \to \infty}{\longrightarrow} \frac{\int_0^\infty \lambda G_0(\theta^*/\lambda) L(d\lambda)}{\int_0^\infty \lambda L(d\lambda)},$$

with $\theta^*$ as before, $L$ is the distribution of popularities, and $G_0$ depends on the inter-arrival distribution $F_0$.

# Universal optimality

## Theorem (In preparation – check ArXiv soon)

Under the above assumptions, the (hard to compute) optimal causal policy converges to a fixed threshold policy with the same limit threshold $\theta^*$.

Therefore, timer policies give a universal asymptotic upper bound on caching/pre-fetching performance.

# Final remarks

- You have to know your traffic before deciding on caching or pre-fetching!

# Final remarks

- You have to know your traffic before deciding on caching or pre-fetching!

- Classical caching is not well suited to regular traffic.

# Final remarks

- You have to know your traffic before deciding on caching or pre-fetching!

- Classical caching is not well suited to regular traffic.

- We identified the hazard rate as the crucial indicator of regularity, and devised a new policy for IHR, that is also asymptotically optimal among all causal policies.

# Final remarks

- You have to know your traffic before deciding on caching or pre-fetching!

- Classical caching is not well suited to regular traffic.

- We identified the hazard rate as the crucial indicator of regularity, and devised a new policy for IHR, that is also asymptotically optimal among all causal policies.

- A lot of open questions, in particular:
    - How we can learn the hazard rates online?
    - How we can estimate the appropriate threshold?
    - What about mixtures of IHR and DHR traffic?

# Thank you!

Andres Ferragut

ferragut@ort.edu.uy

aferragu.github.io

# References I

P. Brémaud.
Point process calculus in time and space.
Springer, NY, 2020.

H. Che, Y. Tung, and Z. Wang.
Hierarchical web caching systems: Modeling, design and experimental results.
IEEE Journal on Selected Areas in Communications, 20(7):1305–1314, 2002.

A. Ferragut, I. Rodriguez, and F. Paganini.
Optimizing TTL caches under heavy tailed demands.
In Proc. of ACM/SIGMETRICS 2016, pages 101–112, June 2016.

A. Ferragut, I. Rodríguez, and F. Paganini.
Optimal timer-based caching policies for general arrival processes.
Queueing Systems, 88(3–4):207–241, 2018.

# References II

📄 N. C. Fofack, P. Nain, G. Neglia, and D. Towsley.
Performance evaluation of hierarchical TTL-based cache networks.
Computer Networks, 65:212–231, 2014.

📄 N. K. Panigrahy, P. Nain, G. Neglia, and D. Towsley.
A new upper bound on cache hit probability for non-anticipative caching policies.
ACM Trans. Model. Perform. Eval. Comput. Syst., 7(2–4), November 2022.