

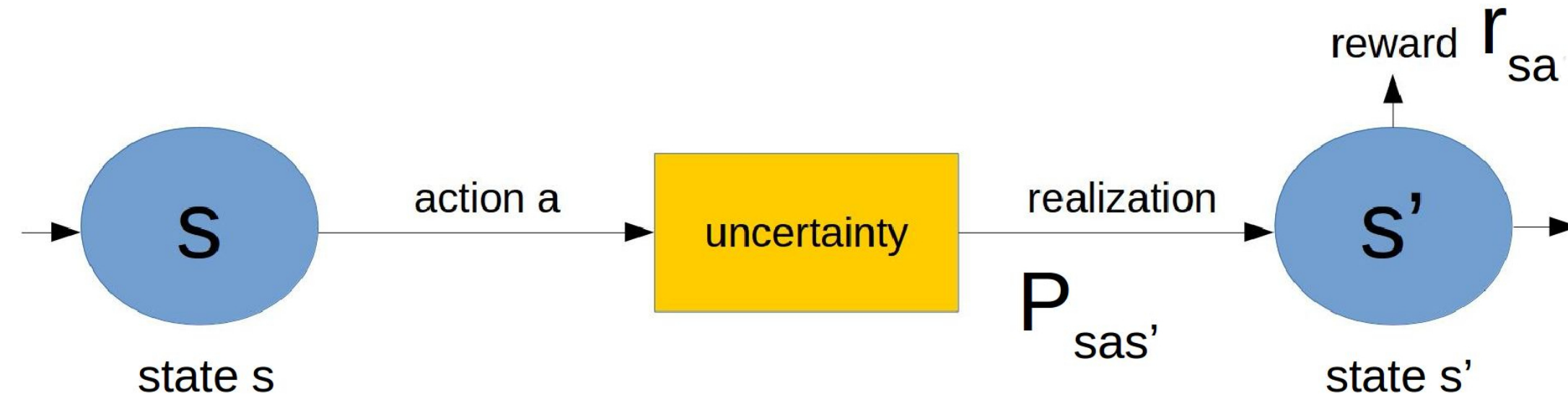
A First-Order Approach to Accelerated Value Iteration

Vineet Goyal, Julien Grand-Clement

{vg2277, jg3728}@columbia.edu, Columbia University

Markov model.

Markov Decision Process: sequential decision making.



Goal: compute optimal sequence of decisions.

There are n states and a discount factor of $\lambda \in (0, 1)$.

Applications: reinforcement learning, healthcare, dynamic pricing, inventory planning, etc.

Classical algorithms.

- Bellman operator $T: \mathbb{R}^n \rightarrow \mathbb{R}^n$, where for $\mathbf{v} \in \mathbb{R}^n$,

$$T(\mathbf{v})_i = \max_a \sum_a \pi_{ia} \cdot (r_{ia} + \lambda \cdot \mathbf{P}_{ia}^\top \mathbf{v}).$$

Goal: find \mathbf{v}^* such that $T(\mathbf{v}^*) = \mathbf{v}^*$.

- Value iteration (VI):

$$\mathbf{v}_0 \in \mathbb{R}^n, \mathbf{v}_{s+1} = T(\mathbf{v}_s), \forall s \geq 0. \quad (\text{VI})$$

→ $\mathbf{v}_s = \mathbf{v}^* + \mathcal{O}(\lambda^s)$ (T is λ -contraction),

→ Problem: running time scales as $\approx 1/(1 - \lambda)$.

Value Iteration and Gradient Descent.

- Gradient Descent:

$$\mathbf{v}_0 \in \mathbb{R}^n, \mathbf{v}_{s+1} = \mathbf{v}_s - \alpha_s \nabla f(\mathbf{v}_s), \forall s \geq 0. \quad (\text{GD})$$

- Idea: $\mathbf{v} - T(\mathbf{v}) = (\mathbf{I} - T)(\mathbf{v}) = \text{gradient of } f: \mathbb{R}^n \rightarrow \mathbb{R}^n$.

- New algorithm:

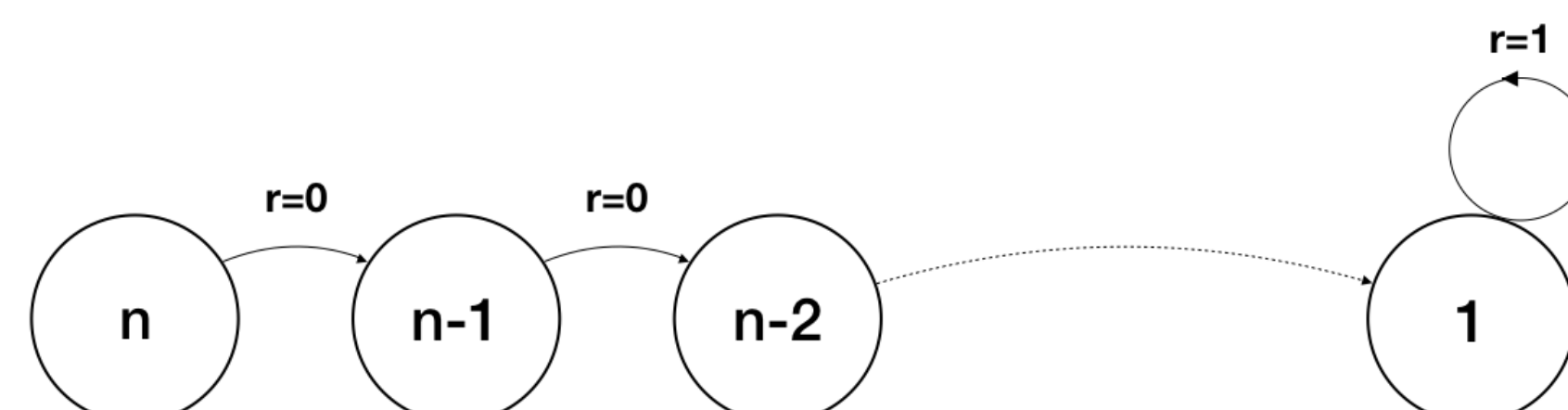
$$\mathbf{v}_0 \in \mathbb{R}^n, \mathbf{v}_{s+1} = \mathbf{v}_s - \alpha_s (\mathbf{v}_s - T(\mathbf{v}_s)), \forall s \geq 0. \quad (\text{R-VI})$$

Recall for μ -strongly convex, L -smooth function f :

$$\mu \leq \frac{\|\nabla f(\mathbf{v}) - \nabla f(\mathbf{w})\|_2}{\|\mathbf{v} - \mathbf{w}\|_2} \leq L. \quad (1)$$

For our Bellman operator T ,

$$(1 - \lambda) \leq \frac{\|(\mathbf{I} - T)(\mathbf{v}) - (\mathbf{I} - T)(\mathbf{w})\|_\infty}{\|\mathbf{v} - \mathbf{w}\|_\infty} \leq (1 + \lambda). \quad (2)$$



Main insights.

- Connection: Convex Optimization and Value Iteration.
- $\mathbf{I} - T$ is the gradient of some function.
- Difficulty: we work with $\|\cdot\|_\infty$ instead of $\|\cdot\|_2$.
- New algorithm: Accelerated Value Iteration (A-VI).

$$\mathbf{v}_0, \mathbf{v}_1 \in \mathbb{R}^n, \begin{cases} \mathbf{h}_s = \mathbf{v}_s + \gamma_s \cdot (\mathbf{v}_s - \mathbf{v}_{s-1}), \\ \mathbf{v}_{s+1} \leftarrow \mathbf{h}_s - \alpha_s (\mathbf{h}_s - T(\mathbf{h}_s)), \end{cases} \quad \forall s \geq 1.$$

Our results.

We tune step sizes using (1) and (2).

- 1 R-VI converges, and running time scales as $\approx 1/(1 - \lambda)$.
- 2 Suppose $T = T^\pi$ (affine operator).
→ If the MDP is *reversible*, A-VI converges and running time scales as $\approx 1/\sqrt{1 - \lambda}$.
→ Without reversibility, A-VI may diverge (ex.: cycle).
- 3 Suppose $T = \text{Bellman operator}$.
→ A-VI is equivalent to Linear Time-Varying Dynamical System.
→ Key role of the joint spectral radius of a set of matrices: a hard problem!

- 4 First-order algorithm:

$$\mathbf{v}_0 = \mathbf{0}, \mathbf{v}_{s+1} \in \text{span}\{\mathbf{v}_0, \dots, \mathbf{v}_s, T(\mathbf{v}_0), \dots, T(\mathbf{v}_s)\}, s \geq 0.$$

Theorem. There exists a hard MDP instance such that no first order algorithm can converge faster than Value Iteration, during the first n iterations.

→ Key difference with convex optimization!

→ Our hard MDP instance:

Numerical experiments.

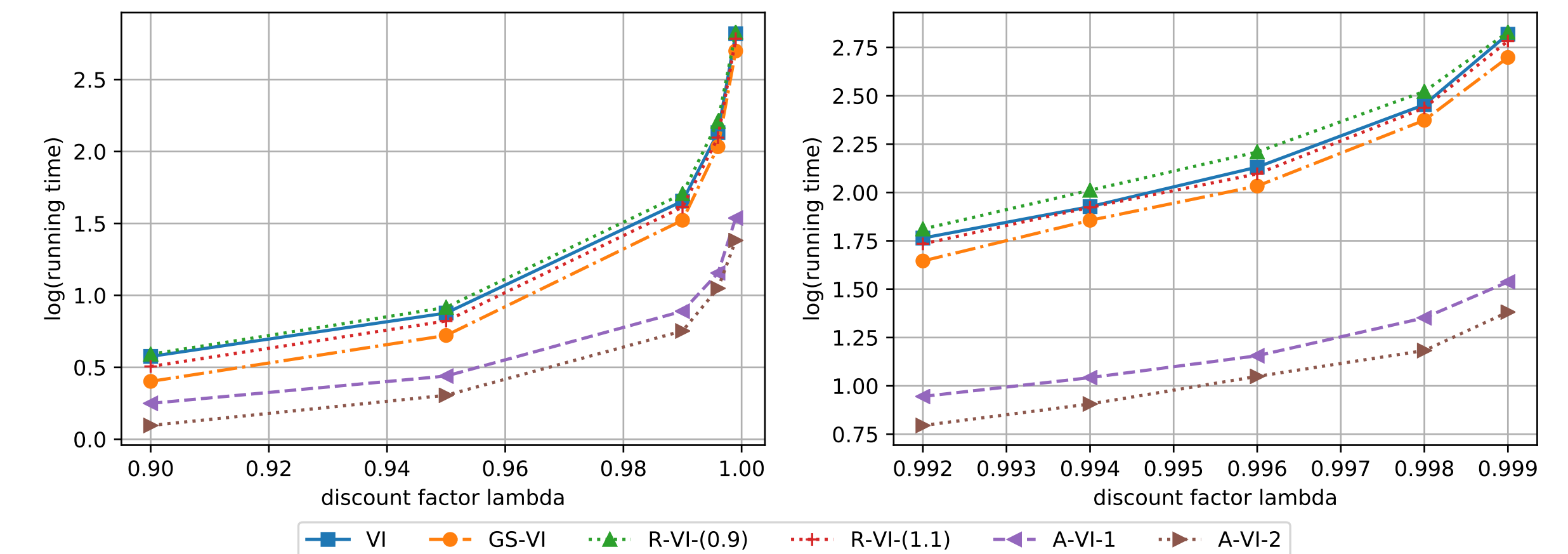


Figure 1: Running time of A-VI vs. state-of-the-art algorithms (log-scale).

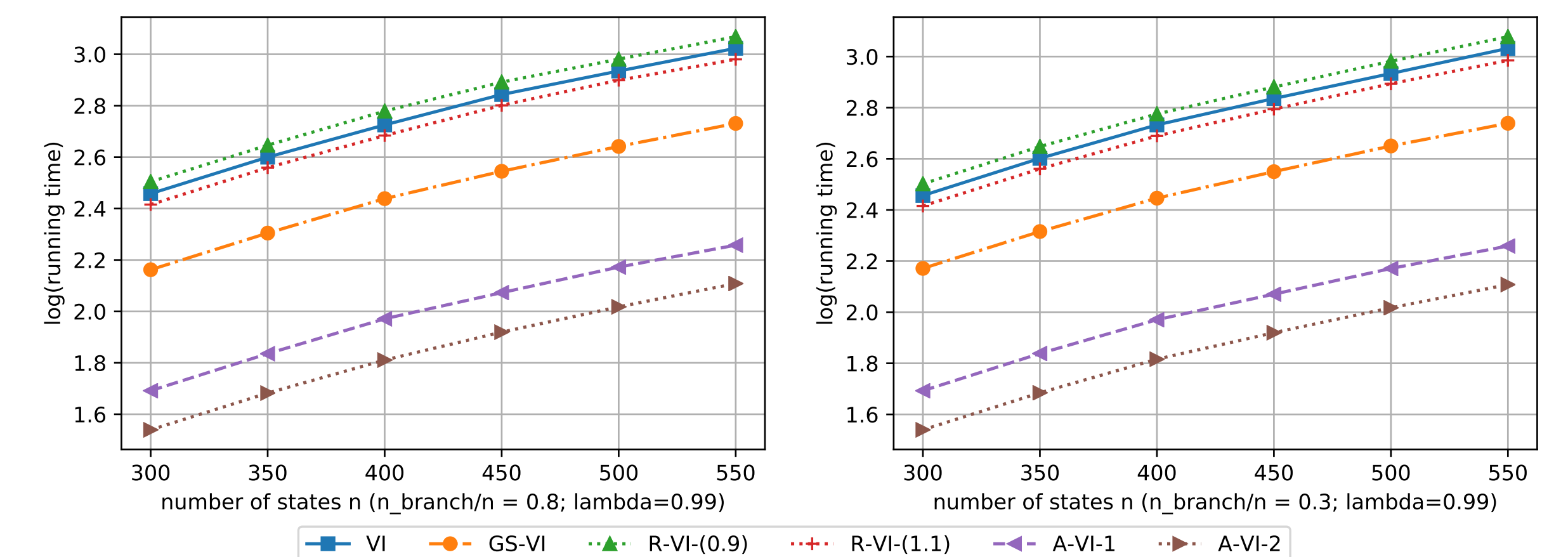


Figure 2: Running time of A-VI and VI on our hard MDP instance.

Reference: *A First-Order Approach to Accelerated Value Iteration*, V. Goyal and J. Grand-Clement, submitted, <https://arxiv.org/abs/1905.09963>.

Webpage: <http://www.columbia.edu/~jg3728/>