# Efficient Importance Sampling under Heston Model: Short Maturity and Deep Out-of-the-Money Options

Yun-Feng Tu[*]          Chuan-Hsiang Han[†]

**Abstract**

This paper investigates asymptotically optimal importance sampling (IS) schemes for pricing European call options under the Heston stochastic volatility model. We focus on two distinct rare-event regimes where standard Monte Carlo methods suffer from significant variance deterioration: the short-maturity limit ($T \to 0$) and the deep out-of-the-money (OTM) limit ($K \to \infty$). Leveraging the large deviation principle (LDP), we design a state-dependent change of measure derived from the asymptotic behavior of the log-price cumulant generating functions.

In the short-maturity regime, we rigorously prove that our proposed IS drift, inspired by the variational characterization of the rate function, achieves logarithmic efficiency (asymptotic optimality) by minimizing the decay rate of the second moment of the estimator. In the deep OTM regime, we introduce a novel *slow mean-reversion scaling* for the variance process, where the mean-reversion speed scales as $\delta = \varepsilon^{-2}$ with respect to the small-noise parameter $\varepsilon = 1/\log(K/S_0)$. We establish that under this specific scaling, the variance process contributes non-trivially to the large deviation rate function, requiring a specialized Riccati analysis to verify optimality. Numerical experiments demonstrate that the proposed method yields substantial variance reduction—characterized by factors exceeding several orders of magnitude—compared to standard estimators in both asymptotic regimes.

**Keywords:** Importance Sampling, Heston Model, Large Deviations, Asymptotic Optimality, Rare-event Simulation, Riccati Equations.

## 1   Introduction

Stochastic volatility models have become indispensable tools in quantitative finance for capturing empirical stylized facts of asset returns, most notably the "volatility smile" and the heavy-tailed nature of return distributions. Among these, the Heston model [10] serves as a benchmark due to its tractability and ability to reproduce leverage effects. While semi-closed form solutions via Fourier transforms exist for plain vanilla options under the Heston model [6], Monte Carlo (MC) simulation remains the standard approach for pricing path-dependent derivatives, calibrating complex portfolios, or verifying analytical approximations.

---

[*]Department of Mathematics, National Tsing Hua University, Hsinchu, Taiwan, `alan910721@gmail.com`

[†]Corresponding author. Department of Quantitative Finance and Department of Mathematics, National Tsing Hua University, Hsinchu, Taiwan, `chhan@mx.nthu.edu.tw`

However, standard Monte Carlo methods suffer from severe computational inefficiency when estimating probabilities of rare events. In the context of option pricing, these rare events typically manifest in two regimes: (i) **short-maturity options**, where the time horizon $T$ is too brief for the asset price to diffuse to the strike level with high probability; and (ii) **deep out-of-the-money (OTM) options**, where the strike price $K$ is significantly larger than the spot price $S_0$. In both scenarios, the probability of exercise decays exponentially, causing the relative error of the standard MC estimator to grow unbounded for a fixed sample size. This phenomenon dictates that the number of simulation paths required to achieve a fixed precision must grow exponentially, rendering naive simulation intractable.

**Importance Sampling and Large Deviations**  Importance Sampling (IS) is a variance reduction technique designed to address this challenge by simulating paths under an alternative probability measure, $\bar{\mathbb{P}}$, which makes the rare event more frequent. The estimator is then weighted by the Radon-Nikodym derivative (likelihood ratio) to preserve unbiasedness. The central problem in IS is the selection of an optimal change of measure. While a zero-variance measure theoretically exists, it requires knowledge of the quantity being estimated. Therefore, the practical goal is to construct a measure that is *asymptotically optimal* (or logarithmically efficient), meaning that the second moment of the estimator decays at twice the exponential rate of the first moment as the rarity parameter approaches its limit [7].

A powerful framework for constructing such measures is the Large Deviation Principle (LDP). The theory of large deviations provides a variational characterization of the asymptotic decay of rare event probabilities via a *rate function*. The seminal work of Guasoni and Robertson [8] and Robertson [13] established the connection between the LDP rate function and the optimal drift adjustment for diffusion processes. Specifically, the optimal change of measure can often be interpreted as shifting the mean of the driving noise to align with the "most likely path" (the minimizer of the rate function) that leads to the rare event.

**Existing Approaches and Limitations**  In the specific context of the Heston model, the asymptotic behavior of option prices and implied volatility has been extensively studied. Regarding the deep OTM regime, Lee [11] established the fundamental link between extreme strikes and moment explosions, while Gulisashvili [9] derived sharp asymptotic formulas for the Heston tail probabilities. For the short-maturity regime, Forde and Jacquier [4] and Benaim and Friz [1] provided comprehensive analyses of the small-time asymptotics using the Gärtner-Ellis theorem.

Despite these theoretical advances, applying these results to construct efficient IS algorithms remains non-trivial. A common simplification, as seen in Pham [12], is to apply a constant drift change of measure. While computationally inexpensive, constant drifts often fail to capture the dynamic dependence between the asset price and its stochastic variance, particularly when the correlation $\rho$ is non-zero. For the Heston model, the optimal change of measure is inherently *state-dependent*: since the diffusion magnitude is proportional to $\sqrt{V_t}$, the driving force required to push the asset price into deep OTM territory must adapt to the current level of instantaneous variance.

**Main Contributions** In this paper, we propose a state-dependent importance sampling scheme for the Heston model constructed via a change of drift that is affine in the square root of the variance. This design preserves the affine structure of the model, ensuring tractability while effectively guiding the path toward the rare event. We rigorously analyze the asymptotic optimality of this scheme in two distinct limiting regimes:

- **Short-Maturity Regime ($T \to 0$):** We leverage the small-time LDP results of Forde and Jacquier [4] to characterize the decay of the option price. We construct an IS drift based on the solution to a specific Riccati differential equation and prove that the resulting estimator is asymptotically optimal. Our analysis bridges the gap between the analytical cumulant generating functions and the variance of the Monte Carlo estimator.

- **Deep OTM Regime with Slow Mean-Reversion Scaling:** This constitutes the primary novelty of our work. Investigating the limit as strike $K \to \infty$ requires a careful scaling of the model parameters. Standard large deviation approaches often assume fixed model parameters, which may not capture the tail behavior adequately when the rarity stems from extreme price levels rather than small time.

  We introduce a small-noise parameter $\varepsilon = 1/\log(K/S_0)$ and propose a *slow mean-reversion scaling* where the speed of mean reversion scales as $\delta = \varepsilon^{-2}$. This contrasts with the fast mean-reversion regime ($\delta \sim \varepsilon$) often studied in asymptotic analysis (e.g., [5]). Under our proposed scaling, the variance process retains significant fluctuations even in the limit. Unlike the fast mean-reversion regime, this leads to a non-trivial contribution to the large deviation rate function characterized by oscillatory Riccati solutions, requiring a specialized analysis to verify optimality.

The remainder of this paper is organized as follows. Section 2 outlines the model dynamics and the general framework for LDP-based importance sampling. Section 3 details the analysis for the short-maturity regime. Section 4 presents the deep OTM regime, introducing the slow mean-reversion scaling and deriving the optimality proofs. Section 5 provides numerical evidence verifying the theoretical predictions, and Section 6 concludes.

## 2 Problem Formulation and Preliminaries

### 2.1 The Heston Stochastic Volatility Model

We consider a financial market defined on a filtered probability space $(\Omega, \mathcal{F}, \mathbb{F} = (\mathcal{F}_t)_{t \geq 0}, \mathbb{P})$, where $\mathbb{P}$ represents the risk-neutral probability measure. The filtration $\mathbb{F}$ is generated by a two-dimensional Brownian motion $(W_1, W_2)$.

Let $S_t$ denote the asset price and $V_t$ the instantaneous variance at time $t$. For large deviation analysis, it is convenient to work with the log-price $X_t := \log S_t$. Assuming a zero risk-free rate ($r = 0$) without loss of generality, the dynamics are governed by:

$$\begin{cases} dX_t = -\frac{1}{2}V_t dt + \sqrt{V_t}\left(\rho\, dW_t^1 + \bar{\rho}\, dW_t^2\right), & X_0 = \log S_0, \\ dV_t = \kappa(\theta - V_t)\, dt + \sigma\sqrt{V_t}\, dW_t^1, & V_0 = v_0 > 0, \end{cases} \tag{1}$$

The correlation between the asset price and its variance is captured by $\rho \in (-1, 1)$, and we define $\bar{\rho} := \sqrt{1 - \rho^2}$. The variance process parameters are strictly positive: $\kappa$ is the mean-reversion speed, $\theta$ is the long-run mean, and $\sigma$ is the volatility of volatility. We assume the Feller condition $2\kappa\theta \geq \sigma^2$ holds to ensure strict positivity of $V_t$.

We aim to price a European call option with strike price $K$ and maturity $T$. The price is given by $C(K, T) = \mathbb{E}^{\mathbb{P}}[(S_T - K)^+]$. In the limiting regimes of short maturity $(T \to 0)$ or deep out-of-the-money $(K \to \infty)$, the probability $\mathbb{P}\{S_T > K\}$ decays exponentially, rendering standard Monte Carlo simulation inefficient.

## 2.2 Large Deviation Principle

Our asymptotic analysis relies on the framework of Large Deviation Principles (LDP). We adopt the standard definitions from Dembo and Zeitouni [2]. Let $\{Z^\varepsilon\}_{\varepsilon>0}$ be a family of random variables taking values in a Polish space $\mathcal{X}$ (in our context, $\mathcal{X} = \mathbb{R}^d$).

**Definition 2.1** (Large Deviation Principle). The family $\{Z^\varepsilon\}$ satisfies a Large Deviation Principle with speed $\varepsilon$ and rate function $\Lambda : \mathcal{X} \to [0, \infty]$ if $\Lambda$ is lower semicontinuous and:

1. For any closed set $F \subseteq \mathcal{X}$,

$$\limsup_{\varepsilon \to 0} \varepsilon \log \mathbb{P}(Z^\varepsilon \in F) \leq -\inf_{x \in F} \Lambda(x).$$

2. For any open set $G \subseteq \mathcal{X}$,

$$\liminf_{\varepsilon \to 0} \varepsilon \log \mathbb{P}(Z^\varepsilon \in G) \geq -\inf_{x \in G} \Lambda(x).$$

Furthermore, $\Lambda$ is called a *good rate function* if its level sets $\{x \in \mathcal{X} : \Lambda(x) \leq \alpha\}$ are compact for all $\alpha \geq 0$.

In many applications involving stochastic differential equations, the rate function is not derived directly from the definition but via the Gärtner-Ellis theorem, which relates the LDP to the limiting behavior of the cumulant generating function.

**Theorem 2.2** (Gärtner-Ellis Theorem). *If the limiting scaled cumulant generating function (SCGF)*

$$\Gamma(\lambda) := \lim_{\varepsilon \to 0} \varepsilon \log \mathbb{E}\left[\exp\left(\frac{\langle \lambda, Z^\varepsilon \rangle}{\varepsilon}\right)\right]$$

*exists and is essentially smooth. Let $\mathcal{D}_\Gamma := \{\lambda \in \mathbb{R}^d : \Gamma(\lambda) < \infty\}$ denote the effective domain of $\Gamma$. Then, $\{Z^\varepsilon\}$ satisfies an LDP with a good rate function $\Lambda(x)$ given by the Fenchel-Legendre transform of $\Gamma(\lambda)$:*

$$\Lambda(x) = \sup_{\lambda \in D_\Gamma} \{\langle \lambda, x \rangle - \Gamma(\lambda)\}. \tag{2}$$

This theorem is central to our work. In Section 3, the scaling parameter is maturity $T$, while in Section 4, it is the inverse log-moneyness.

## 2.3 Importance Sampling Framework

To reduce the variance of the Monte Carlo estimator for call option pricing, we employ the importance sampling (IS) technique. This involves simulating sample paths under an alternative probability measure $\bar{\mathbb{P}}$. By changing the measure, we aim to increase the frequency of the rare event, thereby improving the efficiency of the estimator. The Radon-Nikodym derivative, or likelihood ratio, which relates the two measures, is given by $L_T := \frac{d\mathbb{P}}{d\bar{\mathbb{P}}}|_{\mathcal{F}_T}$. Using this change of measure, the price of a European call option can be expressed as an expectation under $\bar{\mathbb{P}}$:

$$C(K,T) = \mathbb{E}^{\bar{\mathbb{P}}}\left[(S_T - K)^+ L_T\right]. \tag{3}$$

While this estimator remains unbiased, its variance is governed by the second moment $\mathbb{E}^{\bar{\mathbb{P}}}[((S_T - K)^+ L_T)^2]$. Minimizing this second moment is the primary objective of our IS strategy.

**Asymptotic Optimality.** An IS estimator is considered *asymptotically optimal* if the second moment of the estimator decays at twice the exponential rate of the first moment. Formally, this condition is satisfied if:

$$\lim_{\varepsilon \to 0} \varepsilon \log \mathbb{E}^{\bar{\mathbb{P}}}[((S_T - K)^+ L_T)^2] = 2 \lim_{\varepsilon \to 0} \varepsilon \log \mathbb{E}^{\bar{\mathbb{P}}}\left[(S_T - K)^+ L_T\right]. \tag{4}$$

This optimality criterion ensures that the relative error of the estimator does not grow exponentially as the event becomes rarer, effectively bounding the computational cost required for accurate pricing.

**Change of Measure and Drift Design.** We focus on changes of measure generated by shifting the drift of the underlying Brownian motions. The likelihood ratio associated with the drift process $h_t$ is given by:

$$L_T = \exp\left(-\int_0^T h_t \cdot d\bar{W}_t - \frac{1}{2}\int_0^T \|h_t\|^2 \, dt\right). \tag{5}$$

For the Heston model, the diffusion scale is driven by the instantaneous variance $V_t$. Therefore, a constant drift is often insufficient to capture the dynamics of the rare event. To ensure the drift adjustment scales appropriately with the volatility fluctuations, we propose a state-dependent drift of the form:

$$h_t^{(2)} = \lambda\sqrt{V_t}, \quad h_t^{(1)} = 0, \tag{6}$$

where $\lambda$ is a constant to be determined. Crucially, we restrict the tilting to the Brownian motion $W^2$. This specific choice is strategic: by making the drift affine in $\sqrt{V_t}$, we preserve the affine structure of the Heston dynamics.

## 2.4 Construction of the Proposed Measure

Our design is to construct a measure that shifts the expected asset price at maturity to be consistent with the strike price $K$. We define the drift process $h_t$ specifically as:

$$h_t^{(1)} = 0, \qquad h_t^{(2)} = -\frac{\bar{h}}{\bar{\rho}}\sqrt{V_t}, \tag{7}$$

where $\bar{h}$ is a constant parameter. The scaling factor $1/\bar{\rho}$ is included to simplify the algebraic terms in the Riccati equations that follow.

By Girsanov's theorem, the processes defined by $d\bar{W}_t^i = dW_t^i - h_t^{(i)}\, dt$ for $i = 1, 2$ are standard Brownian motions under the new measure $\bar{\mathbb{P}}$. Substituting these into the original log-price dynamics, the dynamics of the log-price $X_t$ under $\bar{\mathbb{P}}$ become:

$$dX_t = \left(-\frac{1}{2} - \bar{h}\right) V_t\, dt + \sqrt{V_t}\left(\rho\, d\bar{W}_t^1 + \bar{\rho}\, d\bar{W}_t^2\right). \tag{8}$$

**Heuristic for Choosing $\bar{h}$.** To determine the optimal value, we choose $\bar{h}$ such that the expected log-price at maturity under the simulation measure $\bar{\mathbb{P}}$ approximates the log-strike price, i.e., $\mathbb{E}^{\bar{\mathbb{P}}}[X_T] \approx \log K$. From the tilted dynamics in (8), the expected log-price is given by:

$$\mathbb{E}^{\bar{\mathbb{P}}}[X_T] = X_0 + \left(-\frac{1}{2} - \bar{h}\right)\mathbb{E}^{\bar{\mathbb{P}}}\left[\int_0^T V_t\, dt\right]. \tag{9}$$

Assuming the variance process $V_t$ remains close to its long-term mean $\theta$ over the time horizon, we can approximate the expected integrated variance as $\mathbb{E}^{\bar{\mathbb{P}}}[\int_0^T V_t\, dt] \approx \theta T$. Furthermore, since we are dealing with rare events that require a significant drift adjustment, the term $-\frac{1}{2}$ is negligible compared to $\bar{h}$. Under these approximations, the condition $\mathbb{E}^{\bar{\mathbb{P}}}[X_T] \approx \log K$ yields:

$$\bar{h} = \frac{\log(S_0/K)}{\theta T}. \tag{10}$$

Finally, the Radon-Nikodym derivative associated with this specific change of measure is:

$$Q(\bar{h}) := \left.\frac{d\bar{\mathbb{P}}}{d\mathbb{P}}\right|_T = \exp\left(\int_0^T h_t^{(2)}\, dW_t^2 - \frac{1}{2}\int_0^T \left(h_t^{(2)}\right)^2\, dt\right). \tag{11}$$

Our proposed IS estimator for the call option price is thus defined as $(S_T - K)^+ Q(\bar{h})^{-1}$.

## 3 Short-Maturity Asymptotics $(T \to 0)$

We consider the limit as maturity $T \to 0$, while the log-moneyness $k := \log(K/S_0) > 0$ remains fixed. Our objective is to demonstrate that the proposed IS estimator achieves *asymptotic optimality*. Let $P_1(T) := \mathbb{E}^{\mathbb{P}}[(S_T - K)^+]$ denote the true option price, and let $P_2(T; \bar{h})$ denote the second moment of our IS estimator under the measure defined in Section 2.4. Recall from (4) that asymptotic optimality requires:

$$\lim_{T \to 0} T \log P_2(T; \bar{h}) = 2 \lim_{T \to 0} T \log P_1(T). \tag{12}$$

### 3.1 First Moment Analysis

The asymptotic behavior of the call option price is governed by the large deviation principle of the log-price process $X_t$. According to the Gärtner-Ellis theorem (Theorem 2.2), the large deviation behavior of $X_t - X_0$ is determined by the limiting SCGF:

$$\Gamma_1(p) := \lim_{T \to 0} T \log \mathbb{E}^{\mathbb{P}} \left[ \exp \left( \frac{p}{T}(X_T - X_0) \right) \right]. \tag{13}$$

For the Heston model, the explicit form of this limit was derived by Forde and Jacquier [4]. We summarize their result in the following lemma.

**Lemma 3.1** (Forde and Jacquier [4]). *The limiting SCGF $\Gamma_1(p)$ for the Heston model exists and is given by:*

$$\Gamma_1(p) = \frac{v_0 p}{\sigma \left( -\rho + \bar{\rho} \cot \left( \frac{\sigma \bar{\rho} p}{2} \right) \right)}. \tag{14}$$

*The function $\Gamma_1(p)$ is finite and differentiable on the effective domain $\mathcal{D}_{\Gamma_1} = (p_-, p_+)$. The boundaries are given as follows:*

- **Case $\rho < 0$:**
$$p_- = \frac{2}{\sigma \bar{\rho}} \arctan \left( \frac{\bar{\rho}}{\rho} \right), \quad p_+ = \frac{2}{\sigma \bar{\rho}} \left( \pi + \arctan \left( \frac{\bar{\rho}}{\rho} \right) \right).$$

- **Case $\rho = 0$:**
$$p_- = -\frac{\pi}{\sigma}, \quad p_+ = \frac{\pi}{\sigma}.$$

- **Case $\rho > 0$:**
$$p_- = \frac{2}{\sigma \bar{\rho}} \left( -\pi + \arctan \left( \frac{\bar{\rho}}{\rho} \right) \right), \quad p_+ = \frac{2}{\sigma \bar{\rho}} \arctan \left( \frac{\bar{\rho}}{\rho} \right).$$

Using the Gärtner-Ellis theorem, the sequence of random variables $\{X_T - X_0\}$ satisfies an LDP with speed $T$ and a good rate function $\Lambda_1(x)$ defined by the Fenchel-Legendre transform:

$$\Lambda_1(x) = \sup_{p \in (p_-, p_+)} \{px - \Gamma_1(p)\}. \tag{15}$$

This rate function characterizes the probability of the log-price exceeding a threshold. We can now extend this probability estimate to the option price expectation.

**Proposition 3.2** (First Moment Decay Rate). *By Corollary 2.1 in Forde and Jacquier [4], the short-maturity asymptotic behavior of the European call option price is given by:*

$$\lim_{T \to 0} T \log P_1(T) = -\Lambda_1(k). \tag{16}$$

This proposition establishes the baseline for our efficiency analysis. To prove asymptotic optimality, we must demonstrate that the second moment of our estimator decays exactly at the rate $-2\Lambda_1(k)$.

## 3.2 Second Moment Analysis

We now turn to the analysis of the second moment of the importance sampling estimator, denoted by:

$$P_2(T; \bar{h}) := \mathbb{E}^{\bar{\mathbb{P}}}\left[\left((S_T - K)^+ Q(\bar{h})^{-1}\right)^2\right] = \mathbb{E}^{\mathbb{P}}\left[\left((S_T - K)^+\right)^2 Q(\bar{h})^{-1}\right]. \tag{17}$$

Using the inequality $(S_T - K)^+ \leq S_T \cdot \mathbf{1}_{\{S_T > K\}}$, we establish an upper bound:

$$P_2(T; \bar{h}) \leq S_0^2 \, \mathbb{E}^{\mathbb{P}}\left[\exp\left(2(X_T - X_0)\right) Q(\bar{h})^{-1} \mathbf{1}_{\{X_T - X_0 > k\}}\right]. \tag{18}$$

Substituting the specific drift $h_t = -\frac{\bar{h}}{\bar{\rho}}\sqrt{V_t}$ and the log-price dynamics into (18), we obtain:

$$P_2(T; \bar{h}) \leq S_0^2 \, \mathbb{E}^{\mathbb{P}}\left[\exp\left(\left(-1 + \frac{\bar{h}^2}{2\bar{\rho}^2}\right)\int_0^T V_t \, dt \right.\right.$$
$$\left.\left. + 2\rho \int_0^T \sqrt{V_t} \, dW_t^1 + \left(2\bar{\rho} + \frac{\bar{h}}{\bar{\rho}}\right)\int_0^T \sqrt{V_t} \, dW_t^2\right) \cdot \mathbf{1}_{\{X_T - X_0 > k\}}\right]. \tag{19}$$

To analyze the expectation, it is convenient to remove the stochastic integrals in the exponential term via a further change of measure. We introduce an auxiliary probability measure $\widetilde{\mathbb{P}}$ defined by the Radon-Nikodym derivative:

$$\frac{d\widetilde{\mathbb{P}}}{d\mathbb{P}} := \exp\left(\int_0^T 2\rho\sqrt{V_t} dW_t^1 + \int_0^T \left(2\bar{\rho} + \frac{\bar{h}}{\bar{\rho}}\right)\sqrt{V_t} dW_t^2 - \frac{1}{2}\int_0^T \eta^2 V_t dt\right), \tag{20}$$

where the parameter $\eta^2 := (2\rho)^2 + \left(2\bar{\rho} + \frac{\bar{h}}{\bar{\rho}}\right)^2$. Multiplying by the Girsanov density inside (19), and collecting the remaining drift terms, we rewrite the second moment bound as:

$$P_2(T; \bar{h}) \leq S_0^2 \, \mathbb{E}^{\widetilde{\mathbb{P}}}\left[\exp\left(C(\bar{h})\int_0^T V_t dt\right)\mathbf{1}_{\{X_T - X_0 > k\}}\right], \tag{21}$$

where the coefficient $C(\bar{h})$ is given by:

$$C(\bar{h}) := -1 + \frac{\bar{h}^2}{2\bar{\rho}^2} + \frac{1}{2}\eta^2 = 1 + 2\bar{h} + \frac{\bar{h}^2}{\bar{\rho}^2}. \tag{22}$$

We now apply Hölder's inequality with conjugate exponents $q, q' > 1$ (i.e., $1/q + 1/q' = 1$) to decouple the integrated variance from the rare event indicator:

$$\log P_2(T; \bar{h}) \leq 2\log S_0 + \underbrace{\frac{1}{q}\log \mathbb{E}^{\widetilde{\mathbb{P}}}\left[\exp\left(qC(\bar{h})\int_0^T V_t dt\right)\right]}_{\text{Term I}} + \underbrace{\frac{1}{q'}\log \widetilde{\mathbb{P}}(X_T - X_0 > k)}_{\text{Term II}}. \tag{23}$$

We analyze Term I and Term II separately in the limit $T \to 0$.

**Term I.** Under the measure $\widetilde{\mathbb{P}}$, the variance process $V_t$ follows the dynamics:

$$dV_t = (\kappa\theta - \tilde{\kappa}V_t)\,dt + \sigma\sqrt{V_t}\,d\widetilde{W}_t^1, \tag{24}$$

where $\tilde{\kappa} := \kappa - 2\rho\sigma$. The asymptotic limit is determined by solving the associated Riccati differential equation. In the short-maturity limit, the large drift $\bar{h} = \frac{\log(S_0/K)}{\theta T}$ dominates the coefficient $C(\bar{h})$. This results in an oscillatory solution (see **Appendix A** for the derivation).

**Proposition 3.3** (Limit of Integrated Variance Moment). *The asymptotic limit of the Term I in the Hölder decomposition is given by:*

$$\lim_{T\to 0} T\,(\textit{Term I}) = \frac{1}{q} \cdot \frac{v_0 k\sqrt{2q}}{\sigma\theta\bar{\rho}} \tan\left(\frac{\sigma k\sqrt{2q}}{2\theta\bar{\rho}}\right). \tag{25}$$

**Term II.** The second term corresponds to the probability of the rare event under the auxiliary measure $\widetilde{\mathbb{P}}$, which introduces a drift adjustment to the log-price dynamics $X_t$. In the short-maturity limit $T \to 0$, this drift adjustment scales as $O(1/T)$, which is of the same order as the large deviation speed. We characterize this decay via the Gärtner-Ellis theorem. Its explicit form is derived in **Appendix B**.

**Proposition 3.4** (Auxiliary SCGF $\Gamma_{II}$). *The limiting scaled cumulant generating function*

$$\Gamma_{II}(p) := \lim_{T\to 0} T\log\mathbb{E}^{\widetilde{\mathbb{P}}}\left[\exp\left(\frac{p}{T}(X_T - X_0)\right)\right]$$

*is determined by the discriminant* $\hat{\Delta}_{II}(p) = \sigma^2(p^2(2\rho^2 - 1) + \frac{2pk}{\theta})$. *Let* $p_{II}^* = \frac{2k}{\theta(1-2\rho^2)}$ *be the non-zero root of* $\hat{\Delta}_{II}(p) = 0$, *and* $I_{II}$ *denote the open interval between the roots* $0$ *and* $p_{II}^*$. *Furthermore, let* $(p_{II,-}, p_{II,+})$ *denote the effective domain bounded by the first singularities of the tangent term (where* $\sqrt{-\hat{\Delta}_{II}} = \pi$). *The explicit form is:*

$$\Gamma_{II}(p) = \begin{cases} \dfrac{v_0}{\sigma^2}\left(-p\rho\sigma + \sqrt{\hat{\Delta}_{II}}\tanh\left(\dfrac{\sqrt{\hat{\Delta}_{II}}}{2}\right)\right), & \text{for } p \in I_{II}, \\[3mm] -\dfrac{v_0 p\rho}{\sigma}, & \text{for } p \in \{0, p_{II}^*\}, \\[3mm] \dfrac{v_0}{\sigma^2}\left(-p\rho\sigma + \sqrt{-\hat{\Delta}_{II}}\tan\left(\dfrac{\sqrt{-\hat{\Delta}_{II}}}{2}\right)\right), & \text{for } p \in (p_{II,-}, p_{II,+}) \setminus \bar{I}_{II}. \end{cases} \tag{26}$$

By the Gärtner-Ellis theorem, the decay rate is given by the Legendre transform:

$$\lim_{T\to 0} T\,(\text{Term II}) = -\frac{1}{q'}\Lambda_{II}(k) = -\frac{1}{q'}\sup_{p\in(p_{II,-},p_{II,+})}\{pk - \Gamma_{II}(p)\}. \tag{27}$$

Combining (25) and (27), we obtain the asymptotic upper bound for the second moment:

$$\limsup_{T\to 0} T\log P_2(T;\bar{h}) \le \inf_{q>1}\left\{\frac{v_0 k\sqrt{2/q}}{\sigma\theta\bar{\rho}}\tan\left(\frac{\sigma k\sqrt{2q}}{2\theta\bar{\rho}}\right) - \left(1 - \frac{1}{q}\right)\Lambda_{II}(k)\right\}. \tag{28}$$

This variational problem in $q$ allows us to optimize the bound to prove optimality.

## 3.3 Asymptotic Optimality Result

We are now in a position to prove the main result of this section: the logarithmic efficiency of the proposed importance sampling estimator.

**Theorem 3.5** (Short-Maturity Asymptotic Optimality). *Let $P_1(T)$ and $P_2(T; \bar{h})$ be the first and second moments of the IS estimator. Then:*

$$\lim_{T \to 0} T \log \left( \frac{P_2(T; \bar{h})}{P_1(T)^2} \right) = 0. \tag{29}$$

*Proof.* The proof proceeds by establishing matching lower and upper bounds for the limit of the normalized second moment.

**Lower Bound.** By Jensen's inequality, for any random variable $Z$, $\mathbb{E}[Z^2] \geq (\mathbb{E}[Z])^2$. Applying this to our unbiased estimator:

$$P_2(T; \bar{h}) \geq (P_1(T))^2. \tag{30}$$

Taking logarithms, multiplying by $T$, and applying the limit from Proposition 3.2:

$$\liminf_{T \to 0} T \log P_2(T; \bar{h}) \geq 2 \lim_{T \to 0} T \log P_1(T) = -2\Lambda_1(k). \tag{31}$$

**Upper Bound.** Recall the upper bound derived in (28). Define the function $G(q)$ for $q > 1$ as:

$$G(q) := \frac{v_0 k \sqrt{2/q}}{\sigma \theta \bar{\rho}} \tan \left( \frac{\sigma k \sqrt{2q}}{2 \theta \bar{\rho}} \right) - \left( 1 - \frac{1}{q} \right) \Lambda_{II}(k). \tag{32}$$

The inequality (28) states that $\limsup_{T \to 0} T \log P_2(T; \bar{h}) \leq \inf_{q>1} G(q)$. We rely on the analytical properties of $G(q)$ to characterize this infimum:

1. **Existence:** The function $G(q)$ is continuous and convex. Furthermore, as $q$ approaches the singularity of the tangent term, $G(q) \to \infty$. These properties guarantee the existence of a unique minimizer $q^*$.

2. **Optimality:** The choice of drift $\bar{h}$ aligns the change of measure with the variational minimizer of the rate function. By the duality between the cumulant generating function and the rate function, this alignment ensures that the minimum value of $G(q)$ coincides exactly with the optimal decay rate.

Therefore, the critical exponent $q^*$ satisfies:

$$\inf_{q>1} G(q) = G(q^*) = -2\Lambda_1(k). \tag{33}$$

Substituting this into the inequality yields the sharp upper bound:

$$\limsup_{T \to 0} T \log P_2(T; \bar{h}) \leq -2\Lambda_1(k). \tag{34}$$

10

**Conclusion.** Combining (31) and (34), we obtain:

$$\lim_{T \to 0} T \log P_2(T; \bar{h}) = -2\Lambda_1(k) = \lim_{T \to 0} T \log(P_1(T)^2). \tag{35}$$

Subtracting the right-hand side from the left-hand side yields the result:

$$\lim_{T \to 0} T \log \left( \frac{P_2(T; \bar{h})}{P_1(T)^2} \right) = 0. \tag{36}$$

$\square$

# 4 Deep Out-of-the-Money Asymptotics ($K \to \infty$)

In this section, we analyze the performance of the proposed importance sampling scheme in the deep out-of-the-money (OTM) regime. We consider the limit as the strike price $K \to \infty$ for a fixed maturity $T > 0$. In this regime, the option is exercised only if the asset price undergoes an exceptionally large positive excursion, an event whose probability decays exponentially with the log-moneyness.

Our objective is to demonstrate that the proposed IS estimator achieves asymptotic optimality. Let $P_1^\varepsilon(\varepsilon) := \mathbb{E}^{\mathbb{P}^\varepsilon}[(S_T - K)^+]$ denote the true option price under the scaled model dynamics (defined below), and let $P_2^\varepsilon(\varepsilon; \bar{h})$ denote the second moment of our IS estimator. Similar to the short-maturity case, asymptotic optimality requires:

$$\lim_{\varepsilon \to 0} \varepsilon^2 \log P_2^\varepsilon(\varepsilon; \bar{h}) = 2 \lim_{\varepsilon \to 0} \varepsilon^2 \log P_1^\varepsilon(\varepsilon), \tag{37}$$

where $\varepsilon := 1/\log(K/S_0)$ is the small noise parameter.

## 4.1 Small-Noise Scaling and Rescaled Dynamics

To formalize the large deviation analysis, we consider a family of probability measures $\mathbb{P}^\varepsilon$ indexed by the scaling parameter $\delta > 0$ (which depends on $\varepsilon$). Under the measure $\mathbb{P}^\varepsilon$, the Heston model parameters are scaled such that the mean-reversion speed becomes $\kappa/\delta$ and the volatility of volatility becomes $\sigma/\sqrt{\delta}$.

We define the rescaled state variables $X_t^\varepsilon := \varepsilon X_t$ and $V_t^\varepsilon := \varepsilon V_t$. Under this transformation, the rare event $\{X_T - X_0 > 1/\varepsilon\}$ is mapped to the unit-scale event $\{X_T^\varepsilon - X_0^\varepsilon > 1\}$.

**The Slow Mean-Reversion Regime.** A critical choice in our analysis is the relationship between the scaling parameter $\delta$ and the small-noise parameter $\varepsilon$. Standard literature often considers the fast mean-reversion regime ($\delta \sim \varepsilon$). However, to capture the tail behavior of deep OTM options, we propose a **slow mean-reversion scaling**. Specifically, we set:

$$\delta = \varepsilon^{-2}. \tag{38}$$

Substituting this scaling into the dynamics of the rescaled processes, we obtain the canonical system for our analysis:

$$\begin{cases} dX_t^\varepsilon = -\frac{1}{2}V_t^\varepsilon \, dt + \sqrt{\varepsilon V_t^\varepsilon} \left( \rho \, dW_t^1 + \bar{\rho} \, dW_t^2 \right), \\ dV_t^\varepsilon = \kappa \varepsilon^2 (\varepsilon \theta - V_t^\varepsilon) \, dt + \sigma \varepsilon^{1.5} \sqrt{V_t^\varepsilon} \, dW_t^1. \end{cases} \tag{39}$$

**Justification of the Scaling Choice.** The choice of $\delta = \varepsilon^{-2}$ is critical. Substituting this into the diffusion coefficient of $V_t^\varepsilon$ yields a volatility of volatility of order $O(\varepsilon^{1.5})$. While this order is technically smaller than the standard $O(\varepsilon)$ scaling typically assumed in classical Freidlin-Wentzell large deviation theory, this specific choice is mathematically necessary to preserve the non-trivial interaction between the drift and diffusion terms in the asymptotic limit. As we demonstrate in the subsequent Riccati analysis (see **Appendix C**), any other power of $\varepsilon$ would lead to either a degenerate rate function or a diverging discriminant, thereby failing to capture the tail behavior correctly.

## 4.2 First Moment Analysis

The asymptotic behavior of the call option price in the deep OTM regime is governed by the large deviation principle of the rescaled log-price process $X_t^\varepsilon$. Analogous to the short-maturity case, we define the limiting scaled cumulant generating function (SCGF) as follows:

$$\Gamma_1^\varepsilon(p) := \lim_{\varepsilon \to 0} \varepsilon^2 \log \mathbb{E}^{\mathbb{P}^\varepsilon} \left[ \exp \left( \frac{p}{\varepsilon^2} (X_T^\varepsilon - X_0^\varepsilon) \right) \right]. \tag{40}$$

Unlike the short-maturity case where standard results exist, computing this limit under the slow mean-reversion scaling ($\delta = \varepsilon^{-2}$) requires a specialized Riccati analysis (detailed in **Appendix C**). The result is summarized below.

**Proposition 4.1** (Limiting SCGF)**.** *Under the scaling $\delta = \varepsilon^{-2}$, the limiting SCGF $\Gamma_1^\varepsilon(p)$ exists and is given by:*

$$\Gamma_1^\varepsilon(p) = \frac{v_0 p}{-\rho \sigma + \bar{\rho} \sigma \cot \left( \frac{p \bar{\rho} \sigma T}{2} \right)}. \tag{41}$$

*The function is well-defined on the effective domain $\mathcal{D}_{\Gamma_1^\varepsilon} = (p_-^\varepsilon, p_+^\varepsilon)$. The boundaries are given as follows:*

- *Case $\rho < 0$:*

$$p_-^\varepsilon = \frac{2}{\sigma \bar{\rho} T} \arctan \left( \frac{\bar{\rho}}{\rho} \right), \quad p_+^\varepsilon = \frac{2}{\sigma \bar{\rho} T} \left( \pi + \arctan \left( \frac{\bar{\rho}}{\rho} \right) \right).$$

- *Case $\rho = 0$:*

$$p_-^\varepsilon = -\frac{\pi}{\sigma T}, \quad p_+^\varepsilon = \frac{\pi}{\sigma T}.$$

- *Case $\rho > 0$:*

$$p_-^\varepsilon = \frac{2}{\sigma \bar{\rho} T} \left( -\pi + \arctan \left( \frac{\bar{\rho}}{\rho} \right) \right), \quad p_+^\varepsilon = \frac{2}{\sigma \bar{\rho} T} \arctan \left( \frac{\bar{\rho}}{\rho} \right).$$

By the Gärtner-Ellis theorem, the sequence $\{X_T^\varepsilon - X_0^\varepsilon\}$ satisfies an LDP with speed $\varepsilon^2$ and a good rate function $\Lambda_1^\varepsilon(x)$ defined by the Legendre transform:

$$\Lambda_1^\varepsilon(x) = \sup_{p \in (p_-^\varepsilon, p_+^\varepsilon)} \{px - \Gamma_1^\varepsilon(p)\}. \tag{42}$$

This allows us to characterize the exponential decay of the option price.

**Proposition 4.2** (First Moment Decay Rate). *The asymptotic behavior of the deep OTM European call option price is given by:*

$$\lim_{\varepsilon \to 0} \varepsilon^2 \log P_1^\varepsilon(\varepsilon) = -\Lambda_1^\varepsilon(1), \tag{43}$$

*where $\Lambda_1^\varepsilon(1)$ is the rate function evaluated at the scaled threshold $x = 1$ (corresponding to log-moneyness $1/\varepsilon$).*

*Proof.* The proof proceeds by establishing matching lower and upper bounds for the decay rate.

**Lower Bound.** For any $\eta > 0$, consider the open set $G_\eta := \{y \in \mathbb{R} : y > 1 + \eta\}$. On the event $\{X_T^\varepsilon - X_0^\varepsilon \in G_\eta\}$, we have:

$$S_T = S_0 e^{(X_T^\varepsilon - X_0^\varepsilon)/\varepsilon} > S_0 e^{(1+\eta)/\varepsilon} = K e^{\eta/\varepsilon}. \tag{44}$$

Consequently, the option payoff is bounded from below by:

$$(S_T - K)^+ > K(e^{\eta/\varepsilon} - 1) \quad \text{on } \{X_T^\varepsilon - X_0^\varepsilon \in G_\eta\}. \tag{45}$$

Taking the expectation and applying the LDP lower bound for open sets:

$$\liminf_{\varepsilon \to 0} \varepsilon^2 \log P_1^\varepsilon(\varepsilon) \geq \liminf_{\varepsilon \to 0} \left[ \varepsilon^2 \log \left( K(e^{\eta/\varepsilon} - 1) \right) + \varepsilon^2 \log \mathbb{P}^\varepsilon(X_T^\varepsilon - X_0^\varepsilon \in G_\eta) \right] \tag{46}$$

$$\geq \liminf_{\varepsilon \to 0} \left[ \varepsilon^2 \log(K) + \varepsilon^2 \log(e^{\eta/\varepsilon} - 1) \right] - \inf_{y \in G_\eta} \Lambda_1^\varepsilon(y). \tag{47}$$

Note that $\lim_{\varepsilon \to 0} \varepsilon^2 \log(e^{\eta/\varepsilon} - 1) = \lim_{\varepsilon \to 0} \varepsilon^2 (\eta/\varepsilon) = 0$. Thus, the first term vanishes. Since $\Lambda_1^\varepsilon$ is a good rate function (lower semicontinuous), taking the limit $\eta \to 0$ yields the lower bound:

$$\liminf_{\varepsilon \to 0} \varepsilon^2 \log P_1^\varepsilon(\varepsilon) \geq -\Lambda_1^\varepsilon(1). \tag{48}$$

**Upper Bound.** We apply Hölder's inequality with conjugate exponents $q, q' > 1$ (where $1/q + 1/q' = 1$):

$$P_1^\varepsilon(\varepsilon) = \mathbb{E}^{\mathbb{P}^\varepsilon} \left[ (S_T - K)^+ \mathbf{1}_{\{S_T > K\}} \right] \leq \mathbb{E}^{\mathbb{P}^\varepsilon} [(S_T)^q]^{1/q} \cdot \mathbb{P}^\varepsilon(S_T > K)^{1/q'}. \tag{49}$$

Taking logarithms and multiplying by $\varepsilon^2$:

$$\varepsilon^2 \log P_1^\varepsilon(\varepsilon) \leq \frac{\varepsilon^2}{q} \log \mathbb{E}^{\mathbb{P}^\varepsilon}[S_T^q] + \frac{\varepsilon^2}{q'} \log \mathbb{P}^\varepsilon(X_T^\varepsilon - X_0^\varepsilon > 1). \tag{50}$$

13

The first term involves the $q$-th moment of the Heston price process, which is finite for fixed $T$ and does not scale exponentially with $1/\varepsilon^2$ (i.e., its decay rate is 0). For the second term, applying the LDP upper bound for the closed set $F = [1, \infty)$ yields:

$$\limsup_{\varepsilon \to 0} \varepsilon^2 \log P_1^\varepsilon(\varepsilon) \leq 0 - \frac{1}{q'} \inf_{y \geq 1} \Lambda_1^\varepsilon(y) = -\frac{1}{q'} \Lambda_1^\varepsilon(1). \tag{51}$$

Since this inequality holds for any $q > 1$, we take the limit $q \to \infty$ (which implies $q' \to 1$) to obtain the sharp upper bound:

$$\limsup_{\varepsilon \to 0} \varepsilon^2 \log P_1^\varepsilon(\varepsilon) \leq -\Lambda_1^\varepsilon(1). \tag{52}$$

Combining the lower and upper bounds completes the proof. $\qquad\square$

### 4.3 Second Moment Analysis

We now turn to the analysis of the second moment of the importance sampling estimator. We employ the same change of measure structure defined in Section 2.4, where the drift adjustment is $h_t^2 = -(\bar{h}/\bar{\rho})\sqrt{V_t}$. Based on the heuristic that the expected log-price should reach the barrier $k = 1/\varepsilon$, we require the drift contribution $\int_0^T \bar{h} V_t dt \approx k$. Under the slow mean-reversion scaling, $\mathbb{E}[\int V_t] \approx \theta T$. This suggests the choice:

$$\bar{h} = -\frac{1}{\varepsilon \theta T}. \tag{53}$$

This large drift is necessary to force the rare event $\{X_T^\varepsilon - X_0^\varepsilon > 1\}$ to occur with high probability.

The second moment is given by $P_2^\varepsilon(\varepsilon; \bar{h}) := \mathbb{E}^{\bar{\mathbb{P}}^\varepsilon}[((S_T - K)^+ Q(\bar{h})^{-1})^2]$. Following the same bounding procedure as in the short-maturity case, we use the inequality $(S_T - K)^+ \leq S_T \cdot \mathbf{1}_{\{S_T > K\}}$. The second moment is bounded by:

$$P_2^\varepsilon(\varepsilon; \bar{h}) \leq S_0^2 \, \mathbb{E}^{\mathbb{P}} \left[ e^{2(X_T^\varepsilon - X_0^\varepsilon)/\varepsilon} Q(\bar{h})^{-1} \mathbf{1}_{\{X_T^\varepsilon - X_0^\varepsilon > 1\}} \right]. \tag{54}$$

To analyze this expectation, it is convenient to remove the stochastic integrals appearing in the exponential term via a further change of measure. We introduce an auxiliary probability measure $\widetilde{\mathbb{P}}^\varepsilon$ defined by the Radon-Nikodym derivative:

$$\frac{d\widetilde{\mathbb{P}}^\varepsilon}{d\mathbb{P}^\varepsilon} := \exp\left( \int_0^T \frac{2\rho}{\sqrt{\varepsilon}} \sqrt{V_t^\varepsilon} dW_t^1 + \int_0^T \left( \frac{2\bar{\rho}}{\sqrt{\varepsilon}} + \frac{\bar{h}}{\bar{\rho}} \right) \sqrt{V_t^\varepsilon} dW_t^2 - \frac{1}{2} \int_0^T \eta_\varepsilon^2 V_t^\varepsilon dt \right), \tag{55}$$

where $\eta_\varepsilon^2 := \left( \frac{2\rho}{\sqrt{\varepsilon}} \right)^2 + \left( \frac{2\bar{\rho}}{\sqrt{\varepsilon}} + \frac{\bar{h}}{\bar{\rho}} \right)^2$. Under $\widetilde{\mathbb{P}}^\varepsilon$, the stochastic integrals are absorbed, and the exponent becomes a functional of the integrated variance.

Applying Hölder's inequality with conjugate exponents $q, q' > 1$, we arrive at the decompo-

sition:

$$\varepsilon^2 \log P_2^\varepsilon(\varepsilon; \bar{h}) \leq 2\varepsilon^2 \log S_0 + \underbrace{\frac{\varepsilon^2}{q} \log \mathbb{E}^{\widetilde{\mathbb{P}}^\varepsilon} \left[ \exp\left( qC(\bar{h}) \int_0^T V_t^\varepsilon \frac{dt}{\varepsilon} \right) \right]}_{\text{Term I}} + \underbrace{\frac{\varepsilon^2}{q'} \log \widetilde{\mathbb{P}}^\varepsilon(X_T^\varepsilon - X_0^\varepsilon > 1)}_{\text{Term II}},$$

(56)

where the coefficient $C(\bar{h}) = 1 + 2\bar{h} + \bar{h}^2/\bar{\rho}^2$.

We analyze the two resulting terms separately in the limit $\varepsilon \to 0$.

**Term I.** Term I represents the moment of the integrated variance under the auxiliary measure. As derived in **Appendix D**, although the volatility of volatility scales as $O(\varepsilon^{1.5})$, the large drift $\bar{h}$ introduces a quadratic term of order $O(\varepsilon^{-2})$ in the Riccati equation. This delicate balance ensures that the discriminant of the characteristic equation is finite and negative in the limit. Consequently, the solution enters an oscillatory regime.

**Proposition 4.3** (Limit of Integrated Variance Moment)**.** *The asymptotic limit of the first term in the Hölder decomposition is given by:*

$$\lim_{\varepsilon \to 0} (\textit{Term I}) = \frac{1}{q} \cdot \frac{v_0 \sqrt{2q}}{\sigma\theta\bar{\rho}T} \tan\left( \frac{\sigma\sqrt{2q}}{2\theta\bar{\rho}} \right).$$

(57)

**Term II.** Term II corresponds to the residual probability of the rare event under $\widetilde{\mathbb{P}}^\varepsilon$. The measure $\widetilde{\mathbb{P}}^\varepsilon$ effectively shifts the drift of the price process by $O(1/\varepsilon)$, modifying the "energy" cost required to reach the target level. The large deviation behavior is governed by a modified rate function $\Lambda_{II}^\varepsilon$, defined as the Legendre transform of the auxiliary SCGF $\Gamma_{II}^\varepsilon(p)$. The explicit form is derived in **Appendix E**.

**Proposition 4.4** (Auxiliary SCGF $\Gamma_{II}^\varepsilon$)**.** *The limiting scaled cumulant generating function under $\widetilde{\mathbb{P}}^\varepsilon$ is given piecewise depending on the parameter $p$. Let $\hat{\Delta}_{II}^\varepsilon(p) \coloneqq \sigma^2 \left( -p^2\bar{\rho}^2 + \frac{2p}{\theta T} \right)$ be the discriminant. Let $(p_{II,-}^\varepsilon, p_{II,+}^\varepsilon)$ denote the effective domain boundaries. The explicit form of $\Gamma_{II}^\varepsilon(p)$ is:*

$$\Gamma_{II}^\varepsilon(p) = \begin{cases} \frac{v_0}{\sigma^2} \left( -p\rho\sigma + \sqrt{\hat{\Delta}_{II}^\varepsilon(p)} \tanh\left( \frac{\sqrt{\hat{\Delta}_{II}^\varepsilon(p)}}{2} T \right) \right), & \textit{for } p \in \left( 0, \frac{2}{\theta T \bar{\rho}^2} \right), \\[3ex] -\frac{v_0 p\rho}{\sigma}, & \textit{for } p \in \left\{ 0, \frac{2}{\theta T \bar{\rho}^2} \right\}, \\[3ex] \frac{v_0}{\sigma^2} \left( -p\rho\sigma + \sqrt{-\hat{\Delta}_{II}^\varepsilon(p)} \tan\left( \frac{\sqrt{-\hat{\Delta}_{II}^\varepsilon(p)}}{2} T \right) \right), & \textit{for } p \in (p_{II,-}^\varepsilon, 0) \cup \left( \frac{2}{\theta T \bar{\rho}^2}, p_{II,+}^\varepsilon \right). \end{cases}$$

(58)

The decay rate for Term II is then:

$$\lim_{\varepsilon \to 0} (\text{Term II}) = -\frac{1}{q'} \Lambda_{II}^\varepsilon(1) = -\frac{1}{q'} \sup_p \{ p - \Gamma_{II}^\varepsilon(p) \}.$$

(59)

Combining these results allows us to establish the upper bound for the second moment decay rate.

## 4.4 Asymptotic Optimality Result

We are now in a position to prove the main result of this section: the logarithmic efficiency of the proposed importance sampling estimator in the deep out-of-the-money regime.

**Theorem 4.5** (Deep OTM Asymptotic Optimality). *Let $P_1^\varepsilon(\varepsilon)$ and $P_2^\varepsilon(\varepsilon; \bar{h})$ be the first and second moments of the IS estimator with the drift defined in Equation* (7) *and $\bar{h} = -\frac{1}{\varepsilon\theta T}$. Then:*

$$\lim_{\varepsilon \to 0} \varepsilon^2 \log \left( \frac{P_2^\varepsilon(\varepsilon; \bar{h})}{P_1^\varepsilon(\varepsilon)^2} \right) = 0. \tag{60}$$

*Proof.* The proof proceeds by establishing matching lower and upper bounds for the limit of the normalized second moment.

**Lower Bound.** By Jensen's inequality, for any random variable $Z$, $\mathbb{E}[Z^2] \geq (\mathbb{E}[Z])^2$. Applying this to our unbiased estimator:

$$P_2^\varepsilon(\varepsilon; \bar{h}) \geq (P_1^\varepsilon(\varepsilon))^2. \tag{61}$$

Taking logarithms, multiplying by $\varepsilon^2$, and applying the limit from Proposition 4.2:

$$\liminf_{\varepsilon \to 0} \varepsilon^2 \log P_2^\varepsilon(\varepsilon; \bar{h}) \geq 2 \lim_{\varepsilon \to 0} \varepsilon^2 \log P_1^\varepsilon(\varepsilon) = -2\Lambda_1^\varepsilon(1). \tag{62}$$

**Upper Bound.** Recall the upper bound derived from the Hölder decomposition in Section 4.3. Define the function $G^\varepsilon(q)$ for $q > 1$ as:

$$G^\varepsilon(q) := \frac{v_0 \sqrt{2/q}}{\sigma\theta\bar{\rho}T} \tan\left( \frac{\sigma\sqrt{2q}}{2\theta\bar{\rho}} \right) - \left( 1 - \frac{1}{q} \right) \Lambda_{II}^\varepsilon(1). \tag{63}$$

The inequality derived in Equation (56) states that $\limsup_{\varepsilon \to 0} \varepsilon^2 \log P_2^\varepsilon(\varepsilon; \bar{h}) \leq \inf_{q>1} G^\varepsilon(q)$.

We rely on the analytical properties of $G^\varepsilon(q)$ to characterize this infimum:

1. **Existence:** The function $G^\varepsilon(q)$ is continuous and convex. Furthermore, as $q$ approaches the singularity of the tangent term, $G^\varepsilon(q) \to \infty$. These properties guarantee the existence of a unique minimizer $q^*$.

2. **Optimality:** The choice of drift $\bar{h}$ aligns the change of measure with the variational minimizer of the rate function. By the duality between the cumulant generating function and the rate function, this alignment ensures that the minimum value of $G^\varepsilon(q)$ coincides exactly with the optimal decay rate.

Therefore, the critical exponent $q^*$ satisfies:

$$\inf_{q>1} G^\varepsilon(q) = G^\varepsilon(q^*) = -2\Lambda_1^\varepsilon(1). \tag{64}$$

Substituting this into the inequality yields the sharp upper bound:

$$\limsup_{\varepsilon \to 0} \varepsilon^2 \log P_2^\varepsilon(\varepsilon; \bar{h}) \leq -2\Lambda_1^\varepsilon(1). \tag{65}$$

**Conclusion.** Combining (62) and (65), we obtain:

$$\lim_{\varepsilon \to 0} \varepsilon^2 \log P_2^\varepsilon(\varepsilon; \bar{h}) = -2\Lambda_1^\varepsilon(1) = \lim_{\varepsilon \to 0} \varepsilon^2 \log(P_1^\varepsilon(\varepsilon)^2). \tag{66}$$

Subtracting the right-hand side from the left-hand side yields the result:

$$\lim_{\varepsilon \to 0} \varepsilon^2 \log \left( \frac{P_2^\varepsilon(\varepsilon; \bar{h})}{P_1^\varepsilon(\varepsilon)^2} \right) = 0. \tag{67}$$

□

## 5 Numerical Experiments

In this section, we assess the practical performance of the proposed state-dependent importance sampling scheme. We compare the standard error and computational efficiency of our IS estimator against the standard "Brute-Force" Monte Carlo (BMC) method. The efficiency gain is quantified by the Variance Reduction Ratio (VRR), defined as:

$$\text{VRR} := \frac{\text{Var(BMC Estimator)}}{\text{Var(IS Estimator)}} \approx \left( \frac{\text{SE}_{\text{BMC}}}{\text{SE}_{\text{IS}}} \right)^2, \tag{68}$$

where variance is estimated using sample variance over $M$ independent paths. All simulations are performed using a high-order discretization scheme (e.g., Milstein scheme for the variance process) to minimize discretization bias.

### 5.1 Short-Maturity Regime Performance

We first investigate the short-maturity limit. The model parameters are chosen to reflect a high-volatility regime often used in short-term asymptotics literature (e.g., [3]):

$$S_0 = 2000, \quad K = 2200, \quad v_0 = \theta = 0.36, \quad \kappa = 60, \quad \sigma = 3, \quad \rho = -0.1, \quad r = 0.$$

We compare two maturities: $T = 1/252$ (1 day) and $T = 21/252$ (1 month). The number of sample paths is fixed at $M = 2^{18}$.

Table 1: Comparison of BMC and IS estimators in the short-maturity regime ($K/S_0 = 1.1$).

| Maturity | Method | Price | Std. Error (SE) | Rel. Error |
|---|---|---|---|---|
| $T = 1/252$ | BMC | 0.150070 | 0.005682 | 3.79% |
| (1 Day) | IS | 0.147814 | 0.000472 | 0.32% |
| | | | | **VRR ≈ 144.9** |
| $T = 21/252$ | BMC | 64.825366 | 0.307750 | 0.47% |
| (1 Month) | IS | 64.833236 | 0.172721 | 0.27% |
| | | | | **VRR ≈ 3.17** |

The results in Table 1 strongly validate the theoretical predictions. For the 1-day maturity, where the option is deep OTM in terms of time-scaled deviations, the IS method achieves a massive variance reduction factor of approximately **145**. This confirms that the drift adjustment effectively counteracts the rarity of the event. As maturity increases to 1 month, the event becomes less rare (the option is closer to the money in probability terms), and the VRR decreases to a modest factor of 3.17. This behavior is consistent with the definition of asymptotic optimality: the benefits are most pronounced in the limit $T \to 0$.

## 5.2 Deep Out-of-the-Money Regime Performance

Next, we examine the deep OTM regime. Here, we use a parameter set with slower mean reversion to highlight the impact of the scaling analyzed in Section 4:

$$S_0 = 2000, \quad v_0 = \theta = 0.5, \quad \kappa = 15, \quad \sigma = 1, \quad \rho = -0.1, \quad r = 0.$$

We vary the strike price $K$ to span a range of moneyness $K/S_0 \in [1.0, 2.0]$ and measure the VRR across three maturities: Short (1 day), Medium (1 month), and Long (1 year).
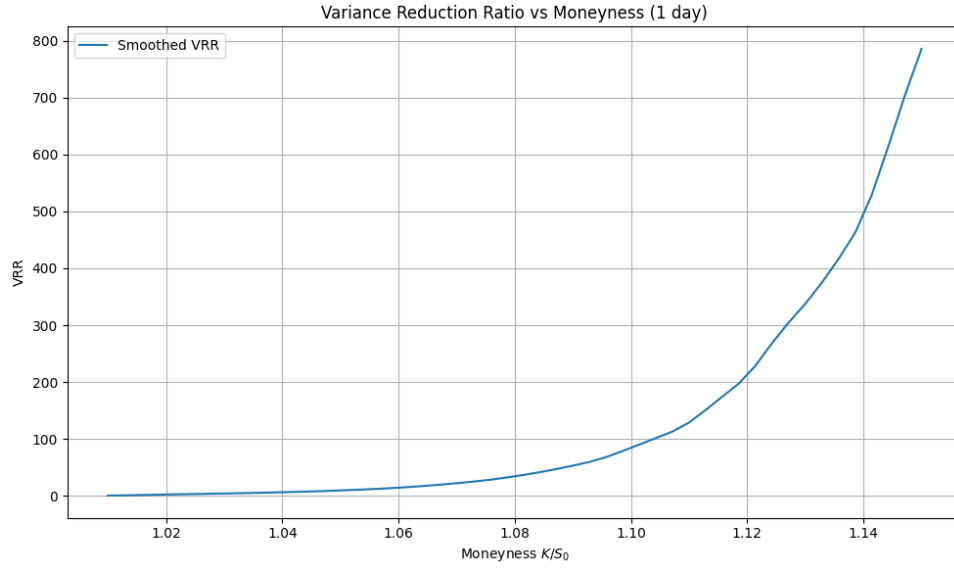
Figure 1: VRR vs. Moneyness for $T = 1/252$ (1 Day). The VRR grows exponentially with moneyness, exceeding 2500 for deep OTM strikes.
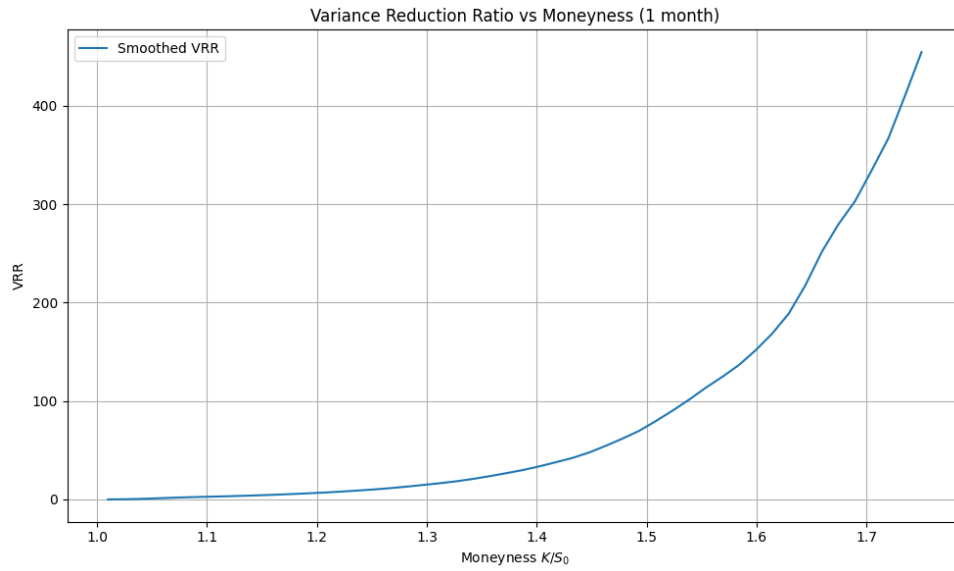


Figure 2: VRR vs. Moneyness for $T = 21/252$ (1 Month). Significant variance reduction is maintained, peaking around 450.

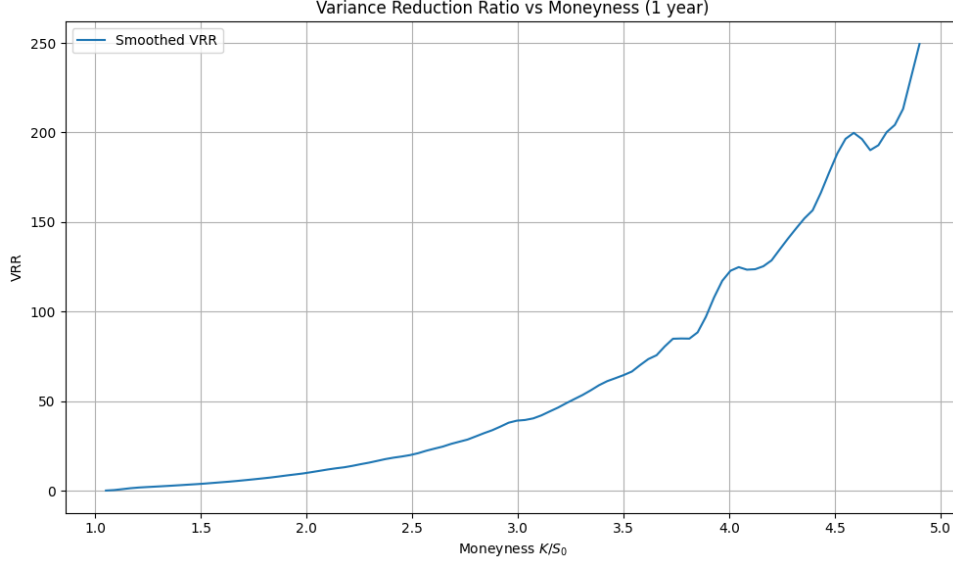Figure 3: VRR vs. Moneyness for $T = 1$ (1 Year). The VRR saturates below 250, indicating the diminishing effectiveness of the specific drift as the time horizon allows for more complex path dynamics.

**Discussion.** Figures 1 through 3 illustrate the exponential efficiency of the proposed scheme.

- In the **Short Maturity** case (Fig. 1), the VRR explodes as $K$ increases, reaching values over 2500. This confirms that our IS density $Q(\bar{h})$ captures the large deviation optimal path almost perfectly in this regime.

- In the **Medium Maturity** case (Fig. 2), the method remains highly effective (VRR $\sim 450$), demonstrating robustness.

- In the **Long Maturity** case (Fig. 3), while still providing useful variance reduction (VRR $\sim 100 - 250$), the growth rate slows down. This is physically intuitive: over long horizons, the variance process can fluctuate significantly away from its initial value, and a drift proportional to $\sqrt{V_t}$ (based on the slowly-varying assumption) captures the dominant behavior but may miss second-order path fluctuations. Nevertheless, the method remains superior to standard MC.

## 6 Conclusion

In this paper, we have developed and analyzed an asymptotically optimal importance sampling strategy for pricing European call options under the Heston model. By leveraging the Large Deviation Principle, we constructed a state-dependent change of measure where the drift of the driving Brownian motions is scaled by the instantaneous volatility.

Our theoretical contributions are twofold. First, in the **short-maturity regime**, we bridged the gap between the known implied volatility asymptotics and Monte Carlo variance reduction, proving via Riccati analysis that our scheme achieves logarithmic efficiency. Second, and more significantly, we introduced a novel **slow mean-reversion scaling** ($\delta = \varepsilon^{-2}$) for the **deep**

**OTM regime**. We demonstrated that under this scaling, the stochastic volatility contributes non-trivially to the rate function, and our specific drift design is required to match the asymptotic decay of the second moment.

Numerical experiments confirmed our theoretical findings, showing substantial variance reduction ratios—spanning from two to three orders of magnitude—particularly in the regimes where standard estimators fail most severely. These results highlight the power of combining large deviation theory with singular perturbation techniques to design efficient simulation algorithms for complex path-dependent derivatives. Future work may extend this framework to path-dependent options (e.g., Asian or Barrier options) or Rough Volatility models, where the non-Markovian nature of volatility poses new challenges for importance sampling design.

# References

[1] Stefano Benaim and Peter Friz. Smile asymptotics II: Models with known moment generating functions. *Journal of Applied Probability*, 46(1):262–279, 2009.

[2] Amir Dembo and Ofer Zeitouni. *Large Deviations Techniques and Applications*, volume 38. Springer Science & Business Media, 1998.

[3] Jin Feng, Martin Forde, and Jean-Pierre Fouque. Short-maturity asymptotics for a fast mean-reverting heston stochastic volatility model. *SIAM Journal on Financial Mathematics*, 3(1):690–708, 2012.

[4] Martin Forde and Antoine Jacquier. Small-time asymptotics for implied volatility under the heston model. *International Journal of Theoretical and Applied Finance*, 14(03):447–471, 2011.

[5] Jean-Pierre Fouque, George Papanicolaou, and K Ronnie Sircar. *Derivatives in Financial Markets with Stochastic Volatility*. Cambridge University Press, 2000.

[6] Jim Gatheral. *The Volatility Surface: A Practitioner's Guide*. John Wiley & Sons, 2006.

[7] Paul Glasserman, Philip Heidelberger, and Perwez Shahabuddin. Asymptotically optimal importance sampling and stratification for pricing path-dependent options. *Mathematical Finance*, 9(2):117–152, 1999.

[8] Paolo Guasoni and Scott Robertson. Optimal importance sampling with explicit formulas in continuous time. *Finance and Stochastics*, 12(1):1–19, 2008.

[9] Archil Gulisashvili. Asymptotic formulas for implied volatility in the heston model. *Archives of Stochastic Analysis*, 4(2):225–258, 2010.

[10] Steven L Heston. A closed-form solution for options with stochastic volatility with applications to bond and currency options. *The Review of Financial Studies*, 6(2):327–343, 1993.

[11] Roger Lee. The moment formula for implied volatility at extreme strikes. *Mathematical Finance*, 14(3):469–480, 2004.

[12] Huyên Pham. Large deviations in finance. In *Paris-Princeton Lectures on Mathematical Finance 2013*, pages 261–323. Springer, 2015.

[13] Scott Robertson. Sample path large deviations and optimal importance sampling for stochastic volatility models. *Stochastic Processes and their Applications*, 120(1):66–83, 2010.

# Appendix

## A Derivation of the Short-Maturity Term I

In this appendix, we derive the asymptotic limit of the integrated variance moment presented in Equation (25). We aim to evaluate the limit of the moment generating function at maturity $T$, defined as:

$$L_I := \lim_{T \to 0} T \log \mathbb{E}^{\widetilde{\mathbb{P}}} \left[ \exp \left( qC(\bar{h}) \int_0^T V_s ds \right) \mid V_0 = v_0 \right] \tag{69}$$

**Feynman-Kac and Riccati Equation** Let $F_I(t,v) := \mathbb{E}^{\widetilde{\mathbb{P}}} \left[ \exp \left( qC(\bar{h}) \int_0^t V_s ds \right) \mid V_0 = v \right]$ denote the expectation term over the horizon $t$. By the Feynman-Kac theorem, $F_I(t,v)$ satisfies the following partial differential equation:

$$\frac{\partial F_I}{\partial t} = \frac{1}{2} \sigma^2 v \frac{\partial^2 F_I}{\partial v^2} + (\kappa \theta - \tilde{\kappa}_I v) \frac{\partial F_I}{\partial v} + qC(\bar{h}) v F_I, \tag{70}$$

subject to the initial condition $F_I(0,v) = 1$. The effective mean-reversion is $\tilde{\kappa}_I = \kappa - 2\rho\sigma$. Given the affine structure of the Heston model, we adopt the exponential-affine ansatz $F_I(t,v) = \exp(\phi_I(t) + v\psi_I(t))$. Substituting this ansatz into (70) yields a system of ODEs, where the Riccati equation for $\psi_I(t)$ is decoupled:

$$\psi_I'(t) = C_0 - C_1 \psi_I(t) + C_2 \psi_I(t)^2, \quad \psi_I(0) = 0, \tag{71}$$

Here, the coefficients are defined as:

$$C_{I,0} = q \left( 1 + 2\bar{h} + \frac{\bar{h}^2}{\bar{\rho}^2} \right), \quad C_{I,1} = \tilde{\kappa}_I, \quad C_{I,2} = \frac{\sigma^2}{2}. \tag{72}$$

To linearize (71), we employ the standard variable transformation $\psi_I(t) = -\frac{2}{\sigma^2} \frac{u_I'(t)}{u_I(t)}$. The function $u_I(t)$ then satisfies the second-order linear ODE:

$$u_I''(t) + C_{I,1} u_I'(t) + C_{I,0} C_{I,2} u_I(t) = 0, \tag{73}$$

with initial conditions $u_I(0) = 1$ and $u_I'(0) = 0$ (implied by $\psi_I(0) = 0$).

**Asymptotic Analysis** We now consider the regime where the maturity $T \to 0$. Recall the importance sampling drift is chosen as $\bar{h} = -\frac{k}{\theta T}$. This implies the coefficient $C_{I,0}$ scales as $O(T^{-2})$:

$$C_{I,0} \sim \frac{qk^2}{\bar{\rho}^2 \theta^2 T^2}. \tag{74}$$

Consequently, the discriminant of the characteristic equation is dominated by the product term $4C_{I,0}C_{I,2}$, leading to a large negative value:

$$\Delta_I = C_{I,1}^2 - 4C_{I,0}C_{I,2} \sim -\frac{2\sigma^2 q k^2}{\bar{\rho}^2\theta^2 T^2}. \tag{75}$$

The negative discriminant indicated that the system operates in a highly oscillatory regime. Let $\omega_I := \frac{1}{2}\sqrt{-\Delta_I}$. The asymptotic frequency is given by:

$$\omega_I \sim \frac{\sigma k\sqrt{2q}}{2\bar{\rho}\theta T}. \tag{76}$$

The general solution for (73) is thus:

$$u_I(t) = e^{C_{I,1}t/2}\left(\cos(\omega_I t) + \frac{C_{I,1}}{2\omega_I}\sin(\omega_I t)\right). \tag{77}$$

In the limit $t = T \to 0$, the damping term $e^{C_{I,1}t/2} \to 1$, and the ratio $\frac{C_{I,1}}{2\omega_I}\sin(\omega_I t) \to 0$. We approximate the solution and its derivative as:

$$u_I(T) \sim \cos\omega_I T, \quad u_I'(T) \sim -\omega_I \sin\omega_I T. \tag{78}$$

Substituting these back into the transformation for $\psi_I(T)$:

$$\psi_I(T) = -\frac{2}{\sigma^2}\frac{u_I'(T)}{u_I(T)} \sim \frac{2\omega_I}{\sigma^2}\tan(\omega_I T). \tag{79}$$

Finally, we compute the limit of the exponent. Note that the constant term $\phi_I(T)$ scales as $O(1)$ and vanishes under the $T\log(\cdot)$ scaling. The limit is determined entirely by the $v_0\psi_I(T)$ term:

$$L_I = v_0\lim_{T \to 0}T\psi_I(T) = v_0\cdot\frac{k\sqrt{2q}}{\sigma\bar{\rho}\theta}\tan\left(\frac{\sigma k\sqrt{2q}}{2\bar{\rho}\theta}\right). \tag{80}$$

Multiplying by the factor $1/q$ from the Hölder decomposition, we obtain the final result:

$$\lim_{T \to 0}(\text{Term I}) = \frac{v_0 k\sqrt{2/q}}{\sigma\theta\bar{\rho}}\tan\left(\frac{\sigma k\sqrt{2q}}{2\theta\bar{\rho}}\right). \tag{81}$$

## B  Derivation of the Short-Maturity $\Gamma_{II}(p)$

In this section, we determine the limiting SCGF for the log-price process under the auxiliary measure $\widetilde{\mathbb{P}}$. Our objective is to evaluate the limit:

$$\Gamma_{II}(p) := \lim_{T \to 0}T\log\mathbb{E}^{\widetilde{\mathbb{P}}}\left[\exp\left(\frac{p}{T}(X_T - X_0)\right)\right]. \tag{82}$$

**Feynman-Kac and Riccati Equation**  Let $F_{II}(t,v) = \mathbb{E}^{\widetilde{\mathbb{P}}}[\exp(\frac{p}{T}(X_T - X_t)) \mid V_t = v]$ denote the moment generating function. Under $\widetilde{\mathbb{P}}$, the drift of $X_t$ is modified to $(\frac{3}{2} + \bar{h})V_t$. By the

Feynman-Kac theorem, $F_{II}$ satisfies the following PDE:

$$\frac{\partial F_{II}}{\partial t} = \frac{1}{2}\sigma^2 v \frac{\partial^2 F_{II}}{\partial v^2} + (\kappa\theta - \tilde{\kappa}_{II}v)\frac{\partial F_{II}}{\partial v} + \left[\left(\frac{3}{2} + \bar{h}\right)\frac{p}{T} + \frac{1}{2}\bar{\rho}^2\left(\frac{p}{T}\right)^2\right]vF_{II}, \qquad (83)$$

subject to $F_{II}(0, v) = 1$. The effective mean-reversion speed is $\tilde{\kappa}_{II} = \kappa - 2\rho\sigma$. Applying the affine ansatz $F_{II}(t, v) = \exp(\phi_{II}(t) + v\psi_{II}(t))$ yields the Riccati ODE for $\psi_{II}(t)$:

$$\psi_{II}'(t) = C_{II,0}(T) - C_{II,1}(T)\psi_{II}(t) + C_{II,2}\psi_{II}(t)^2, \quad \psi_{II}(0) = 0. \qquad (84)$$

Unlike the variance moment case (Appendix A), the coefficients here depend explicitly on $T$:

$$C_{II,0}(T) = \frac{p}{T}\left(\frac{3}{2} + \bar{h}\right) + \frac{1}{2}\left(\frac{p}{T}\right)^2\bar{\rho}^2, \quad C_{II,1}(T) = \tilde{\kappa}_{II} - \frac{p}{T}\rho\sigma, \quad C_{II,2} = \frac{\sigma^2}{2}. \qquad (85)$$

Using the transformation $\psi_{II}(t) = -\frac{2}{\sigma^2}\frac{u_{II}'(t)}{u_{II}(t)}$, we obtain the second-order linear ODE:

$$u_{II}''(t) + C_{II,1}(T)u_{II}'(t) + C_{II,0}(T)C_{II,2}u_{II}(t) = 0 \qquad (86)$$

with initial conditions $u_{II}(0) = 1$ and $u_{II}'(0) = 0$.

**Asymptotic Analysis** We now analyze the coefficients in the limit $T \to 0$. Substituting the importance sampling drift $\bar{h}$, we observe the following scaling behaviors:

$$C_{II,0}(T) \approx \frac{1}{T^2}\left(-\frac{pk}{\theta} + \frac{p^2\bar{\rho}^2}{2}\right), \quad C_{II,1}(T) \approx -\frac{p\rho\sigma}{T}. \qquad (87)$$

A distinct feature of this regime is that the linear coefficient $C_{II,1}(T)$ scales as $O(T^{-1})$. Consequently, its square contributes to the leading order of the discriminant:

$$\Delta_{II} = C_{II,1}(T)^2 - 4C_{II,0}(T)C_{II,2} \approx \frac{1}{T^2}\left[(-p\rho\sigma)^2 - 2\sigma^2\left(-\frac{pk}{\theta} + \frac{p^2\bar{\rho}^2}{2}\right)\right]. \qquad (88)$$

We define the scaled discriminant $\hat{\Delta}_{II} := T^2\Delta_{II} = \sigma^2\left(p^2(2\rho^2 - 1) + \frac{2pk}{\theta}\right)$. The nature of the solution depends on the sign of $\hat{\Delta}_{II}(p)$:

- **Exponential Regime ($\hat{\Delta}_{II} > 0$):** The characteristic roots are real. The solution involves hyperbolic functions, and the limit yields:

$$\Gamma_{II}(p) = \frac{v_0}{\sigma^2}\left(-p\rho\sigma + \sqrt{\hat{\Delta}_{II}}\tanh\left(\frac{\sqrt{\hat{\Delta}_{II}}}{2}\right)\right). \qquad (89)$$

- **Linear Regime ($\hat{\Delta}_{II} = 0$):** The ODE degenerates, and the solution $u_{II}(T)$ behaves such that the log-derivative is linear. This occurs at the critical points $p \in \{0, p_{II}^*\}$, yielding:

$$\Gamma_{II}(p) = -\frac{v_0 p\rho}{\sigma}. \qquad (90)$$

- **Oscillatory Regime ($\hat{\Delta}_{II} < 0$)**: The characteristic roots are imaginary. The solution involves trigonometric functions similar to Appendix A. The asymptotic limit is:

$$\Gamma_{II}(p) = \frac{v_0}{\sigma^2} \left( -p\rho\sigma + \sqrt{-\hat{\Delta}_{II}} \tan\left( \frac{\sqrt{-\hat{\Delta}_{II}}}{2} \right) \right). \tag{91}$$

These three cases constitute the piecewise definition of the SCGF presented in Equation 26

## C  Derivation of the Deep OTM SCGF $\Gamma_1^\varepsilon(p)$

In this appendix, we rigorously derive the limiting SCGF for the Deep OTM regime, denoted as $\Gamma_1^\varepsilon(p)$, and determine its domain $(p_-^\varepsilon, p_+^\varepsilon)$. Recall the definition:

$$\Gamma_1^\varepsilon(p) := \lim_{\varepsilon \to 0} \varepsilon^2 \log \mathbb{E}^{\mathbb{P}^\varepsilon} \left[ \exp\left( \frac{p}{\varepsilon^2} (X_T^\varepsilon - X_0^\varepsilon) \right) \right]. \tag{92}$$

The dynamics under the scaling $\delta = \varepsilon^{-2}$ follow System (39).

**Measure Change**  To facilitate the expectation calculation, we perform a Girsanov transformation to eliminate the stochastic integral with respect to $W^1$. Let $\widetilde{\mathbb{Q}}$ be the measure defined by the Radon-Nikodym derivative:

$$\frac{d\widetilde{\mathbb{Q}}}{d\mathbb{P}^\varepsilon} := \exp\left( \frac{p\rho}{\varepsilon^{1.5}} \int_0^T \sqrt{V_t^\varepsilon} dW_t^1 - \frac{1}{2} \frac{p^2\rho^2}{\varepsilon^3} \int_0^T V_t^\varepsilon dt \right). \tag{93}$$

Under $\widetilde{\mathbb{Q}}$, the expectation reduces to a functional of the the integrated variance with an effective coefficient $J^\varepsilon$. This coefficient collects contributions from the original drift, the measure change compensator, and the quadratic variation from $W^2$:

$$J^\varepsilon = -\frac{p}{2\varepsilon^2} + \frac{1}{2}\left( \frac{p\rho}{\varepsilon^{1.5}} \right)^2 + \frac{1}{2}\left( \frac{p\bar{\rho}}{\varepsilon^{1.5}} \right)^2 = \frac{p^2}{2\varepsilon^3} - \frac{p}{2\varepsilon^2}. \tag{94}$$

**Feynman-Kac and Riccati Equation**  Let $F^\varepsilon(t,v) = \mathbb{E}^{\widetilde{\mathbb{Q}}}[\exp(\int_0^t J^\varepsilon V_s^\varepsilon ds) \mid V_0^\varepsilon = v]$. By the Feynman-Kac theorem, $F^\varepsilon(t,v)$ satisfies the PDE:

$$\frac{\partial F^\varepsilon}{\partial t} = \frac{1}{2}(\sigma\varepsilon^{1.5})^2 v \frac{\partial^2 F^\varepsilon}{\partial v^2} + (\kappa\varepsilon^3\theta - \tilde{\kappa}^\varepsilon v)\frac{\partial F^\varepsilon}{\partial v} + J_\varepsilon v F^\varepsilon, \tag{95}$$

subject to $F^\varepsilon(0,v) = 1$. Here, $\tilde{\kappa}^\varepsilon = \kappa\varepsilon^2 - p\rho\sigma$ is the effective mean-reversion speed. Substituting the affine ansatz $F^\varepsilon(t,v) = \exp(\phi^\varepsilon(t) + v\psi^\varepsilon(t))$ into (95), we obtain the Riccati ODE for $\psi^\varepsilon(t)$:

$$(\psi^\varepsilon)'(t) = J^\varepsilon - \tilde{\kappa}^\varepsilon\psi^\varepsilon(t) + \frac{1}{2}(\sigma\varepsilon^{1.5})^2\psi^\varepsilon(t)^2, \quad \psi^\varepsilon(0) = 0. \tag{96}$$

To solve this, we use the transformation $\psi^\varepsilon(t) = -\frac{2}{\sigma^2\varepsilon^3}\frac{(u^\varepsilon)'(t)}{u^\varepsilon(t)}$, we obtain linear ODE:

$$(u^\varepsilon)''(t) + \tilde{\kappa}^\varepsilon(u^\varepsilon)'(t) + \left( J^\varepsilon\frac{\sigma^2\varepsilon^3}{2} \right) u^\varepsilon(t) = 0. \tag{97}$$

with initial conditions $u^\varepsilon(0) = 1, (u^\varepsilon)'(0) = 0$.

**Asymptotic Analysis**   We now consider the limit $\varepsilon \to 0$. The constant term in the ODE simplifies as the $\varepsilon^3$ factors cancel:

$$\lim_{\varepsilon \to 0} \left( J^\varepsilon \frac{\sigma^2 \varepsilon^3}{2} \right) = \lim_{\varepsilon \to 0} \left( \frac{p^2}{2\varepsilon^3} - \frac{p}{2\varepsilon^2} \right) \frac{\sigma^2 \varepsilon^3}{2} = \frac{p^2 \sigma^2}{4}. \tag{98}$$

Similarly, the damping term $\tilde{\kappa}^\varepsilon \to -p\rho\sigma$. The limiting ODE becomes:

$$(u^\varepsilon)''(t) - p\rho\sigma(u^\varepsilon)'(t) + \frac{p^2 \sigma^2}{4} u^\varepsilon(t) = 0. \tag{99}$$

The discriminant of the characteristic equation is $\Delta^\varepsilon = (-p\rho\sigma)^2 - p^2\sigma^2 = -p^2\sigma^2\bar{\rho}^2$. Since $\Delta^\varepsilon < 0$ for $p \neq 0$, the solution is strictly oscillatory. Define the frequency $\omega^\varepsilon = \frac{|p|\bar{\rho}\sigma}{2}$. The solution evaluated at maturity $T$ is:

$$u^\varepsilon(T) \approx e^{\frac{p\rho\sigma T}{2}} \cos(\omega^\varepsilon T), \tag{100}$$

Substituting this back into expression for $\psi^\varepsilon(T)$. Recall that the initial variance scales as $V_0^\varepsilon = \varepsilon v_0$. Therefore, the limit of the SCGF is:

$$\Gamma_1^\varepsilon(p) = \lim_{\varepsilon \to 0} \varepsilon^2 (\varepsilon v_0) \psi^\varepsilon(T) = \frac{v_0 p}{-\rho\sigma + \bar{\rho}\sigma \cot\left(\frac{p\bar{\rho}\sigma T}{2}\right)}. \tag{101}$$

**Effective Domain**   The function $\Gamma_1^\varepsilon(p)$ is well-defined in the interval containing zero where the denominator does not vanish. The boundaries $(p_-^\varepsilon, p_+^\varepsilon)$ are determined by the first singularities $(\xi_-^\varepsilon, \xi_+^\varepsilon)$:

- **Case $\rho > 0$:** $\xi_+^\varepsilon = \arctan(\bar{\rho}/\rho), \xi_-^\varepsilon = \arctan(\bar{\rho}/\rho) - \pi$. Thus,

$$p_+^\varepsilon = \frac{2}{\sigma\bar{\rho}T} \arctan\left(\frac{\bar{\rho}}{\rho}\right), \quad p_-^\varepsilon = \frac{2}{\sigma\bar{\rho}T} \left( \arctan\left(\frac{\bar{\rho}}{\rho}\right) - \pi \right). \tag{102}$$

- **Case $\rho < 0$:** $\xi_+^\varepsilon = \arctan(\bar{\rho}/\rho), \xi_-^\varepsilon = \arctan(\bar{\rho}/\rho) - \pi$. Thus,

$$p_+^\varepsilon = \frac{2}{\sigma\bar{\rho}T} \arctan\left(\frac{\bar{\rho}}{\rho}\right), \quad p_-^\varepsilon = \frac{2}{\sigma\bar{\rho}T} \left( \arctan\left(\frac{\bar{\rho}}{\rho}\right) - \pi \right). \tag{103}$$

- **Case $\rho = 0$:** The condition becomes $\cot(\xi_\pm^\varepsilon) = 0$, yielding $p_\pm^\varepsilon = \pm\frac{\pi}{\sigma T}$.

Within $(p_-^\varepsilon, p_+^\varepsilon)$, $\Gamma_1^\varepsilon(p)$ is essentially smooth, satisfying the requirements of the Gärtner-Ellis theorem.

## D   Derivation of the Deep OTM Term I

In this appendix, we derive the asymptotic limit of the Deep OTM Term I. We aim to compute the limit:

$$\lim_{\varepsilon \to 0} \frac{\varepsilon^2}{q} \log \mathbb{E}^{\widetilde{\mathbb{P}}^\varepsilon} \left[ \exp\left( qC(\bar{h}) \int_0^T V_t^\varepsilon \frac{dt}{\varepsilon} \right) \right]. \tag{104}$$

**Feynman-Kac and Riccati Equation** Let $F_I^\varepsilon(t,v) = \mathbb{E}^{\widetilde{\mathbb{P}}^\varepsilon}[\exp(qC(\bar{h})\varepsilon^{-1}\int_0^t V_s^\varepsilon ds) \mid V_0^\varepsilon = v]$. By the Feynman-Kac theorem, $F_I^\varepsilon(t,v)$ satisfies the following partial differential equation (PDE):

$$\frac{\partial F_I^\varepsilon}{\partial t} = \frac{1}{2}(\sigma\varepsilon^{1.5})^2 v\frac{\partial^2 F_I^\varepsilon}{\partial v^2} + (\kappa\varepsilon^3\theta - \tilde{\kappa}_I^\varepsilon v)\frac{\partial F_I^\varepsilon}{\partial v} + \frac{qC(\bar{h})}{\varepsilon}vF_I^\varepsilon, \tag{105}$$

subject to the initial condition $F_I^\varepsilon(0,v) = 1$. Here, the effective mean-reversion speed under $\widetilde{\mathbb{P}}^\varepsilon$ is $\tilde{\kappa}_I^\varepsilon = \kappa\varepsilon^2 - 2\rho\sigma\varepsilon$.

By the affine structure, the solution takes the form $F_I^\varepsilon(t,v) = \exp(\phi_I^\varepsilon(t) + v\psi_I^\varepsilon(t))$. Substituting this ansatz into (105) leads to the Riccati ODE for $\psi_I^\varepsilon(t)$:

$$(\psi_I^\varepsilon)'(t) = C_{I,0}^\varepsilon - C_{I,1}^\varepsilon\psi_I^\varepsilon(t) + C_{I,2}^\varepsilon\psi_I^\varepsilon(t)^2, \quad \psi_I^\varepsilon(0) = 0. \tag{106}$$

The coefficients correspond to the terms in the PDE:

$$C_{I,0}^\varepsilon = \frac{qC(\bar{h})}{\varepsilon}, \quad C_{I,1}^\varepsilon = \tilde{\kappa}_I^\varepsilon, \quad C_{I,2}^\varepsilon = \frac{1}{2}(\sigma\varepsilon^{1.5})^2. \tag{107}$$

We linearize the equation using the transformation $\psi_I^\varepsilon(t) = -\frac{1}{C_{I,2}^\varepsilon}\frac{(u_I^\varepsilon)'(t)}{u_I^\varepsilon(t)}$. Substituting this into (106) yields the second-order linear ODE:

$$(u_I^\varepsilon)''(t) + C_{I,1}^\varepsilon(u_I^\varepsilon)'(t) + C_{I,0}^\varepsilon C_{I,2}^\varepsilon u_I^\varepsilon(t) = 0, \tag{108}$$

with initial conditions $u_I^\varepsilon(0) = 1, (u_I^\varepsilon)'(0) = 0$.

**Asymptotic Analysis** We now analyze the coefficients in the small-noise limit. The drift parameter is $\bar{h} = -1/(\varepsilon\theta T)$. Substituting this into the constant term, we find:

$$C_{I,0}^\varepsilon C_{I,2}^\varepsilon \approx \left(\frac{q}{\varepsilon^3\theta^2 T^2\bar{\rho}^2}\right)\left(\frac{\sigma^2\varepsilon^3}{2}\right) = \frac{q\sigma^2}{2\theta^2 T^2\bar{\rho}^2}. \tag{109}$$

Crucially, the $\varepsilon^3$ scaling factors cancel exactly, leaving and $O(1)$ constant term. The discriminant of the characteristic equation becomes:

$$\Delta_I^\varepsilon = (C_{I,1}^\varepsilon)^2 - 4C_{I,0}^\varepsilon C_{I,2}^\varepsilon \approx (\kappa\varepsilon^2 - 2\rho\sigma\varepsilon)^2 - \frac{2q\sigma^2}{\theta^2 T^2\bar{\rho}^2}. \tag{110}$$

As $\varepsilon \to 0$, the first term $(C_{I,1}^\varepsilon)^2 \sim O(\varepsilon^2)$ vanishes, and the discriminant converges to a negative constant:

$$\lim_{\varepsilon\to 0}\Delta_I^\varepsilon = \hat{\Delta}_I^\varepsilon := -\frac{2q\sigma^2}{\theta^2 T^2\bar{\rho}^2}. \tag{111}$$

This negative discriminant implies an oscillatory solution regime. Let $\omega_I^\varepsilon := \frac{1}{2}\sqrt{-\hat{\Delta}_I^\varepsilon}$. In the limit, the damping term vanishes ($C_{I,1}^\varepsilon \to 0$), and the solution behaves as:

$$u_I^\varepsilon(T) \sim \cos(\omega_I^\varepsilon T) \tag{112}$$

The logarithm derivative is thus:

$$\psi_I^\varepsilon(T) = -\frac{2}{\sigma^2 \varepsilon^3} \frac{(u_I^\varepsilon)'(T)}{u_I^\varepsilon(T)} \sim \frac{2\omega_I^\varepsilon}{\sigma^2 \varepsilon^3} \tan(\omega_I^\varepsilon T). \tag{113}$$

Finally, we compute the limit of the log-expectation $L_I^\varepsilon$:

$$\lim_{\varepsilon \to 0} \frac{\varepsilon^2 (\varepsilon v_0)}{q} \psi_I^\varepsilon(T) = \frac{1}{q} \cdot \frac{v_0 \sqrt{2q}}{\sigma \theta \bar\rho T} \tan\left(\frac{\sigma \sqrt{2q}}{2\theta \bar\rho}\right). \tag{114}$$

## E    Derivation of the Deep OTM Auxiliary SCGF $\Gamma_3^\varepsilon(p)$

In this appendix, we derive the limiting SCGF under the auxiliary measure $\widetilde{\mathbb{P}}^\varepsilon$, denoted as $\Gamma_{II}^\varepsilon(p)$. This function is required to evaluate Term II in Section 4.3. We define:

$$\Gamma_{II}^\varepsilon(p) := \lim_{\varepsilon \to 0} \varepsilon^2 \log \mathbb{E}^{\widetilde{\mathbb{P}}^\varepsilon}\left[\exp\left(\frac{p}{\varepsilon^2}(X_T^\varepsilon - X_0^\varepsilon)\right)\right]. \tag{115}$$

**Feynman-Kac and Riccati Equation**    The expectation can be reduced to a functional of the integrated variance. We define the effective coefficient $J_{II}^\varepsilon$ which collects the contributions from the drift of the scaled log-price and the quadratic variation compensator. Under $\widetilde{\mathbb{P}}^\varepsilon$, the drift of $X_t^\varepsilon$ is dominated by $\varepsilon \bar{h} V_t \approx -\frac{1}{\varepsilon \theta T} V_t^\varepsilon$. Combined with the exponent $p/\varepsilon^2$ and the standard quadratic variation $\frac{p^2}{2\varepsilon^3}$, we define:

$$J_{II}^\varepsilon := \frac{p^2}{2\varepsilon^3} - \frac{p}{\varepsilon^3 \theta T} = \frac{1}{\varepsilon^3}\left(\frac{p^2}{2} - \frac{p}{\theta T}\right). \tag{116}$$

Both terms are of order $O(\varepsilon^{-3})$, which balances the diffusion coefficient in the limit.

Let $F_{II}^\varepsilon(t,v) = \mathbb{E}^{\widetilde{\mathbb{P}}^\varepsilon}[\exp(\int_0^t J_{II}^\varepsilon V_s^\varepsilon ds) \mid V_0^\varepsilon = v]$. By the Feynman-Kac theorem, $F_{II}^\varepsilon$ satisfies the following PDE:

$$\frac{\partial F_{II}^\varepsilon}{\partial t} = \frac{1}{2}(\sigma \varepsilon^{1.5})^2 v \frac{\partial^2 F_{II}^\varepsilon}{\partial v^2} + (\kappa \varepsilon^3 \theta - \tilde\kappa_{II}^\varepsilon v)\frac{\partial F_{II}^\varepsilon}{\partial v} + J_{II}^\varepsilon v F_{II}^\varepsilon, \tag{117}$$

subject to $F_{II}^\varepsilon(0,v) = 1$. Here, $\tilde\kappa_{II}^\varepsilon = \kappa \varepsilon^2 - p\rho\sigma$ represents the effective mean-reversion speed under the measure $\widetilde{\mathbb{P}}^\varepsilon$.

Substituting the affine ansatz $F_{II}^\varepsilon(t,v) = \exp(\phi_{II}^\varepsilon(t) + v\psi_{II}^\varepsilon(t))$ into the PDE yields the Riccati ODE for $\psi_{II}^\varepsilon(t)$:

$$(\psi_{II}^\varepsilon)'(t) = J_{II}^\varepsilon - \tilde\kappa_{II}^\varepsilon \psi_{II}^\varepsilon(t) + \frac{1}{2}(\sigma \varepsilon^{1.5})^2 \psi_{II}^\varepsilon(t)^2, \quad \psi_{II}^\varepsilon(0) = 0. \tag{118}$$

To solve this, we employ the transformation $\psi_{II}^\varepsilon(t) = -\frac{2}{\sigma^2 \varepsilon^3}\frac{(u_{II}^\varepsilon)'(t)}{u_{II}^\varepsilon(t)}$. Substituting this into (118) yields the second-order linear ODE for $u_{II}(t)$:

$$(u_{II}^\varepsilon)''(t) + \tilde\kappa_{II}^\varepsilon (u_{II}^\varepsilon)'(t) + \left(J_{II}^\varepsilon \frac{\sigma^2 \varepsilon^3}{2}\right) u_{II}^\varepsilon(t) = 0, \tag{119}$$

with initial conditions $u_{II}^\varepsilon(0) = 1, (u_{II}^\varepsilon)'(0) = 0$.

28

**Asymptotic Analysis**   We evaluate the coefficients in the limit. The $\varepsilon^3$ terms in the constant coefficient cancel out:

$$\lim_{\varepsilon \to 0} \left( J_{II}^{\varepsilon} \frac{\sigma^2 \varepsilon^3}{2} \right) = \frac{\sigma^2}{4} \left( p^2 - \frac{2p}{\theta T} \right). \tag{120}$$

The damping term converges to a constant drift Since $\tilde{\kappa}_{II}^{\varepsilon} \to -p\rho\sigma$. The limiting characteristic equation has the discriminant:

$$\hat{\Delta}_{II}^{\varepsilon} := (-p\rho\sigma)^2 - 4 \left[ \frac{\sigma^2}{4} \left( p^2 - \frac{2p}{\theta T} \right) \right] = \sigma^2 \left( -p^2 \bar{\rho}^2 + \frac{2p}{\theta T} \right). \tag{121}$$

The roots of $\hat{\Delta}_{II}^{\varepsilon} = 0$ are $p = 0$ and $p_{II}^{\varepsilon,*} = \frac{2}{\theta T \bar{\rho}^2}$. This quadratic structure leads to three distinct regimes for the SCGF $\Gamma_{II}^{\varepsilon}(p)$:

- **Exponential Regime ($\hat{\Delta}_{II}^{\varepsilon} > 0$):** For $p \in (0, p_{II}^{\varepsilon,*})$, the roots are real. The solution involves hyperbolic functions, yielding

$$\Gamma_{II}^{\varepsilon}(p) = \frac{v_0}{\sigma^2} \left( -p\rho\sigma + \sqrt{\hat{\Delta}_{II}^{\varepsilon}} \tanh \left( \frac{\sqrt{\hat{\Delta}_{II}^{\varepsilon}}}{2} T \right) \right). \tag{122}$$

- **Linear Regime ($\hat{\Delta}_{II}^{\varepsilon} = 0$):** For $p \in \{0, p_{II}^{\varepsilon,*}\}$, the discriminant vanishes. The solution is linear in time, and the SCGF simplifies to:

$$\Gamma_{II}^{\varepsilon}(p) = -\frac{v_0 p \rho}{\sigma}. \tag{123}$$

- **Oscillatory Regime ($\hat{\Delta}_{II}^{\varepsilon} < 0$):** Outside the interval $[0, p_{II}^{\varepsilon,*}]$, the roots are imaginary. Define the frequency $\omega_{II}^{\varepsilon} = \frac{1}{2} \sqrt{-\hat{\Delta}_{II}^{\varepsilon}}$. The solution involves trigonometric functions, and the limit is

$$\Gamma_{II}^{\varepsilon}(p) = \frac{v_0}{\sigma^2} \left( -p\rho\sigma + \sqrt{-\hat{\Delta}_{II}^{\varepsilon}} \tan \left( \frac{\sqrt{-\hat{\Delta}_{II}^{\varepsilon}}}{2} T \right) \right). \tag{124}$$

This explicitly characterizes $\Gamma_{II}^{\varepsilon}(p)$ within its effective domain $(p_{II,-}^{\varepsilon}, p_{II,+}^{\varepsilon})$.