

The Black Hole Idealization Paradox

Dominic Ryder¹

¹*Department of Philosophy, Logic and Scientific Method, London School of Economics. d.ryder@lse.ac.uk*

April 18, 2024

Abstract

Stephen Hawking's derivation of Hawking radiation relied on one particular spacetime model, that of a star collapsing into a black hole which then remains in existence forever. He then argued that Hawking radiation implies this model should be thrown away in favour of a different model, that of an evaporating black hole. This aspect of Hawking's argument is an example of an idealization that is pervasive in the literature on black hole thermodynamics, but which has not yet been widely discussed by philosophers. The aim of this paper is to clarify the nature of Hawking's idealization, and to show a sense in which it leads to a paradox. After identifying this idealization paradox in classic derivations of Hawking radiation, I go on to show how various research programmes in black hole thermodynamics can be viewed as possible resolutions to the paradox. I give an initial analysis of the prospects for success of these various resolutions, and show how they shed light on both the philosophical foundations of both Hawking radiation on the nature of idealizations in physics.

1 Introduction

“He must, so to speak, throw away the ladder after he has climbed up it.” – Ludwig Wittgenstein, *Tractatus Logico-Philosophicus*

Derivations of Hawking radiation are cornerstones of modern physics. The consensus view is that Hawking radiation leads to the black hole information paradox, and huge amounts of work in physics has been dedicated to understanding and resolving it (Page, 1994; Raju, 2022). Philosophers have also analysed various aspects of Hawking radiation, including: the black hole information paradox (Belot, Earman, & Ruetsche, 1999; Maudlin, 2017; Manchak & Weatherall, 2018; Wallace, 2020); black hole thermodynamics (Dougherty & Callender, 2016; Wüthrich, 2019; Wallace, 2018, 2019; Prunkl & Timpson, 2019); and the universality of Hawking radiation (Gryb et al., 2019).

In this paper I will argue that there is another problem, distinct from those listed above, that arises because a seemingly essential idealization is used in three mainstream derivations of Hawking radiation: Hawking’s original derivation (1975)¹, Fredenhagen and Haag’s (1990) “watertight” derivation, and algebraic approaches such as Dimock and Kay (1987) and Dappiaggi et al. (2011). This paper establishes the paradox for these derivations, categorises its possible resolutions, and offers an initial analysis of the success of various resolutions. The resolution of this problem, which I call the *idealization paradox*, can teach us about the kinds of idealizations used in science, how global spacetime structure encodes local spacetime structure, and the nature of Hawking radiation.

The paradox arises out of an argument of Hawking (1975), who derived the eponymous radiation in a spacetime which represents a star that collapses into a black hole which, once formed, is unchanging and exists for the rest of time. I will call this spacetime *collapse-Schwarzschild*. In the same paper, Hawking also presented the first arguments that the backreaction of the radiation on the spacetime will lead to a negative energy flux into the black hole, thus causing the black hole to lose mass and evaporate. Given that a black hole evaporates, Hawking reasoned, it is not well represented by collapse-Schwarzschild. Instead, we should represent the black hole using a spacetime that models an evaporating black hole. I will call this spacetime *evaporation-Schwarzschild*. In other words, the use of collapse-Schwarzschild in the derivation of Hawking radiation is an idealization. Significantly, evaporation-Schwarzschild does not exhibit the same properties as collapse-Schwarzschild, and as I show in section 3.2, Hawking’s derivation cannot be carried out in evaporation-Schwarzschild. We throw away the spacetime we were using as a ladder to Hawking radiation, collapse-Schwarzschild, in favour of evaporation-Schwarzschild, but the original derivation is not consistent with our new spacetime.

The thesis of this paper is that, according to Hawking’s and other mainstream derivations, Hawking radiation is inconsistent with black hole evaporation. It is possible to state a sketch of the paradox (the details of which I will complete in the next sections) for a derivation of Hawking radiation based

¹And *a fortiori* Wald (1975), as this is just a more mathematically rigorous version of Hawking’s derivation.

upon a certain set of properties X :

The Idealization Paradox

1. **(Hawking Radiation Derivation)** If spacetime exhibits the set of properties X , then Hawking radiation occurs.
2. **(Backreaction Arguments)** If Hawking radiation occurs, then black hole evaporation occurs.
3. **(Inconsistency Claim)** If black hole evaporation occurs, then spacetime does not exhibit the set of properties X .
4. **(Spacetime Postulate)** Spacetime exhibits the set of properties X .

This set of premises is inconsistent. What justifies the first premise? Of Hawking’s original calculations, Unruh (2014) writes they are “mathematically unimpeachable”, and the other derivations discussed in this paper only improve upon the degree of mathematical rigour. Thus, the secure mathematical status of the derivations in question means the first premise is hard to challenge.² What about the second premise? Using global definitions of energy one can derive a positive energy flux out toward infinity in the black hole spacetimes, as I discuss further in section 2.2. Hence, assuming global energy conservation, one recovers a negative energy flux into the black hole, which is strong motivation for black hole evaporation. The third premise, Inconsistency Claim, is defended in the bulk of this paper. So what about the fourth premise? It is hard to reject Spacetime Postulate, because then we can’t use the derivation of the first premise to derive Hawking radiation. So without the fourth premise, we lose our motivation for believing in Hawking radiation. Hence, according to the derivation used in the first premise, Hawking radiation is inconsistent with evaporation.

Notice that Hawking Radiation Derivation makes reference to a particular derivation.³ I call derivations to which the idealization paradox applies *evaporation-inconsistent derivations*, and conversely those to which it does not *evaporation-consistent derivations*. Thus, I argue Hawking (1975), Fredenhagen and Haag (1990) and algebraic approaches such as Dimock and Kay (1987), Dappiaggi et al. (2011) are evaporation-inconsistent.

Due to the possibility of evaporation-consistent derivations (and resolutions to the paradox for evaporation-inconsistent derivations), the idealization paradox does not imply we ought to be skeptical about the existence of Hawking radiation or black hole evaporation. To protect these phenomena from the paradox, one may claim that there exists a derivation of Hawking radiation that uses physically reasonable properties and is evaporation-consistent. Call this existence claim the *consistency conjecture*. The ‘physically reasonable’ qualification is necessary because a physically implausible evaporation-consistent derivation (e.g. a derivation in a two-dimensional spacetime) should not al-

²Unruh also calls the calculations “nonsense physically” due to the *trans-Planckian problem*, recently discussed by philosophers Gryb et al. (2019). However, I ignore the trans-Planckian problem for the purpose of this paper.

³More precisely, a particular set of properties assumed in a derivation. If two derivations of Hawking radiation assume the exact same properties, then we consider them equivalent for the purpose of this paper.

leviate our concerns. I expect the consistency conjecture is true.⁴ However, I will argue that even if the consistency conjecture is true, the idealization paradox is still paradoxical and must be resolved; the paradox identifies a mystery about how and why evaporation-inconsistent derivations were successful.

Some philosophers, such as Batterman (2002, 2005, 2011, 2017) and Morrison (2012), have argued that idealizations (construed as false descriptions) are essential for our scientific theories and models to represent and explain reality: there is “something deeply correct about the “unrealistic” idealization” (Batterman, 2005, p. 237). Conversely, many have defended the view that idealizations are dispensable (for example, Norton (2012); Butterfield (2011); Menon and Callender (2013) and Palacios (2019, 2022)). This attitude is captured in what Jones (2006)⁵ has called *Earman’s principle*: “no effect can be counted as a genuine physical effect if it disappears when the idealizations are removed” (Earman, 2004, p. 191).⁶

Applying Earman’s principle to Hawking radiation, dispensabilists will presumably demand that the collapse-Schwarzschild idealization must be removed.⁷ This would have the further benefit of helping to explain why new derivations of Hawking radiation continue to be produced, despite Hawking’s original derivation being widely viewed as successfully establishing the phenomenon. However, as we shall soon see, deidealizing Hawking’s argument is not conceptually straightforward, lending some initial plausibility to the essentialist claim. Nonetheless, a more careful look at recent research programmes in the foundations of Hawking radiation also reveals several distinct options for the dispensabilist.

My aim in this paper will be to establish the paradox for the three derivations and then categorise possible dispensabilist responses to the paradox, each of which seeks to deidealize the derivations. The plausible resolutions are associated with prominent research programmes in the foundations of Hawking radiation, including an appeal to quantum gravity, the approximation regime proposed by Hawking (1975), and what I call “essential structure” derivations. I give an initial analysis of these approaches and find their prospects of success vary significantly. In particular Hawking’s approximation regime fails for his own derivation, but essential structure derivations represent a very

⁴In the literature there are derivations which are plausible candidates for evaporation-consistency (e.g. Visser (2003); Parikh and Wilczek (2000)). I discuss these in section 6, but a full analysis requires another paper, which the author intends to undertake in the future of the project. See Curiel (2023) for an overview of the plethora of Hawking radiation derivations.

⁵See also Landsman (2013); Fletcher (2020).

⁶A widely discussed example in this literature is the unrealistic use of infinite limits in statistical mechanics to recover singularities in the thermodynamic theory of phase transitions. For a topical introduction to the debate and bibliography see Shech (2018, 2023). Fletcher et al. (2019); Shech (2018, 2023) catalogue some of the philosophical issues that arise from the use of idealizations in physics; and see Frigg and Hartmann (2020); Potochnik (2017) and Frigg (2022, chapter 11) for general overviews on idealization in science.

⁷Earman’s principle has also been applied to Hawking radiation by Gryb et al. (2019), in which the authors note that the response to the trans-Planckian problem which models Hawking radiation as Goldstone bosons has only been carried out in stationary spacetimes, and it is an important open question whether these models can be deidealized.

promising possible resolution. The lessons of the paradox vary across possible resolutions, but initial hints suggest insight into: the nature of Hawking radiation, how global physical properties encode local physical properties, and what sort of idealizations are used in our best physical theories.

In section 2 I introduce the theory of black holes and quantum field theory on black hole spacetimes that will be required. In section 3 I show that the idealization paradox applies to Hawking’s derivation, and in sections 4 and 5 I show that the idealization paradox bites for Fredenhagen and Haag’s derivation and algebraic approaches respectively. Finally, in section 6 I categorise and analyse resolutions to the paradox.

2 Primer on Quantum Field Theory in Black Hole Spacetimes

This section reviews the background material important to the claim that the derivations I analyse are evaporation-inconsistent. I begin with black hole physics treated from the global perspective (Hawking and Ellis (1973) and Wald (1984)). I then introduce quantum field theory on curved spacetimes which underwrites the particle concept in Hawking radiation, before sketching the black hole evaporation heuristic. Readers familiar with quantum field theory on curved spacetime may wish to skip to section 3.

2.1 Black Hole Spacetimes and Conformal Diagrams

For my purposes, a black hole spacetime is one that is asymptotically flat at past and future null infinity (\mathcal{I}^\pm) and for which there is a region of the spacetime causally isolated from the rest of the spacetime for all time ($J^-(\mathcal{I}^+) \neq M$).⁸ The black hole region is the causally isolated region ($\mathcal{B} = M - J^-(\mathcal{I}^+)$), and the event horizon bounds this region ($\mathcal{H}^E = \partial\mathcal{B}$). There are in fact multiple inequivalent ways to define a black hole (Curiel, 2019), but the one given here is standard in the formal and philosophical foundations of general relativity. Birkhoff’s theorem states that any solution of Einstein’s vacuum ($R_{ab} = T_{ab} = 0$) equations which is spherically symmetric⁹ in an open set V , is isometric in V to part of the inextendible Schwarzschild solution (Hawking & Ellis, 1973). This solution describes an uncharged, non-rotating black hole of mass m with the event horizon at the Schwarzschild radius, $r = 2m$. A spacetime is said to be *inextendible* if there does not exist a ‘larger’ spacetime into which there is a proper isometric embedding.

A spacetime is *stationary* if it admits a global timelike Killing vector field.¹⁰ Roughly, a stationary

⁸Heuristically, a spacetime is asymptotically flat at \mathcal{I}^\pm iff it is approximately Minkowski at infinity, approaches Minkowski smoothly, and is complete. For a formal definition see Wald (1984, chapter 11.1), and see Landsman (2021, sec. 10.3) for in depth discussion of these conditions.

⁹Admits the group $SO(3)$ as a group of isometries, with the group orbits spacelike two-surfaces.

¹⁰A Killing vector field is a vector field whose flow is a one-parameter group of isometries ϕ_t ($\phi_t : M \rightarrow M$ such that $\phi_t^*g = g$).

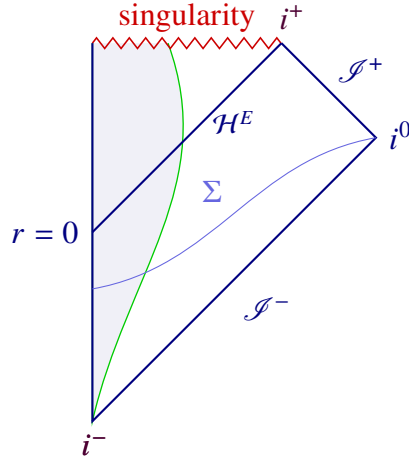


Figure 1: The conformal diagram for stellar collapse into a Schwarzschild black hole. The shaded region represents matter undergoing collapse.

spacetime does not change if one follows the integral curves of the Killing vector field. Schwarzschild spacetime is stationary. However, physical black holes are formed by some astrophysical process, such as stellar collapse, so the physical spacetime will be neither stationary nor vacuum. Therefore, for a more realistic representation, we analyse a spacetime that includes spherically symmetric, non-rotating and uncharged matter that collapses into a black hole. The collapse-Schwarzschild conformal diagram is the resulting model, depicted in figure 1. This diagram will be important, so I identify some of its distinctive features. Outside the matter the spacetime is isometric to Schwarzschild by Birkhoff's theorem, and thus is stationary. Inside the matter the metric will be complicated and non-stationary. The spacetime is globally hyperbolic, meaning that it admits a Cauchy surface and thus a well-posed initial value description.¹¹ One such Cauchy surface is denoted Σ in figure 1. Given any foliation into Cauchy surfaces, once one surface has intersected the event horizon all subsequent surfaces also will. Thus, these models describe black holes which exist forever after their formation.

The final version of a Schwarzschild black hole to consider is evaporation-Schwarzschild. However, first we need Hawking radiation.

2.2 Quantum Fields in Black Hole Spacetimes

We now turn to how particles are defined in quantum field theory, and how non-stationary spacetimes lead to particle creation. This is the core of Hawking radiation according to the mainstream view. The idea is to quantise a classical field theory by defining a Hilbert space, \mathcal{H} , with respect to a time translation symmetry, giving a particle interpretation of the field. Because time translation

¹¹A Cauchy surface is one such that every causal curve (without an endpoint) intersects it exactly once. Therefore, heuristically, a Cauchy surface registers some information about every point in spacetime, and a globally hyperbolic spacetime is causally well-behaved. See also Geroch (1970).

symmetries are generally local in curved spacetimes, the particle interpretation is generally different in different regions, and this leads to particle creation.

In more detail, one begins by modelling a massless complex-valued scalar field, Φ , obeying the covariant wave-equation: $g_{ab}\nabla^a\nabla^b\Phi = 0$. We can now take any of a variety of paths to define a quantum field theory, but roughly one defines a Hilbert space by selecting a subset of the solutions to the covariant wave-equation to represent physical solutions.¹² In stationary spacetimes there is a preferred non-arbitrary way to select this subspace. Namely, we can define a global time coordinate associated with a Killing vector field that characterises time translation symmetry, and choose \mathcal{H} to be the space of positive frequency solutions with respect to this time coordinate (exactly as in Minkowski spacetime for an inertial time coordinate). By non-arbitrarily fixing \mathcal{H} , we non-arbitrarily fix a particle interpretation for our QFT.¹³ Thus, there is a preferred, global definition of a particle in stationary spacetimes. However, in general curved spacetimes there will not be a time translation symmetry which we can exploit to define positive frequency solutions. Therefore, there will not exist a non-arbitrary way to define \mathcal{H} ; so there will be no unique, global definition of a particle. This applies to collapse-Schwarzschild, which is non-stationary.

The central idea of Hawking radiation is that the failure of a spacetime to yield a global preferred particle interpretation leads to particle creation. There are local Killing vector fields at past infinity and future infinity but these differ due to non-stationarity in the bulk region of the spacetime. Therefore, given a vacuum state in the past one has particle content in the future. This is the basis of the Hawking (1975) derivation of Hawking radiation. The details are saved for section 3.1, but in summary: we define \mathcal{H}^\pm on \mathcal{I}^\pm and then choose Φ such that the state is vacuum on \mathcal{I}^- ; by calculating the unitary operator $U : \mathcal{H}^- \rightarrow \mathcal{H}^+$, we can determine the particle number for Φ on \mathcal{I}^+ with respect to \mathcal{H}^+ , and one finds that there is particle creation. Specifically, a thermal spectrum of particles is found at \mathcal{I}^+ . This leads to our next topic, evaporation.¹⁴

Black hole evaporation cannot be directly inferred from the claim that black holes radiate, because they do not radiate like normal black bodies: no part of Hawking radiation lies in the causal future of the black hole. Instead, evaporation is thought to occur due to the backreaction of the radiation and,

¹²For example, following Wald (1995, p.38), first define a state space for a quantum theory called a ‘‘one-particle structure’’. The covariant wave equation admits a symplectic vector space of complex-valued solutions, $(\mathcal{S}^\mathbb{C}, \Omega)$, where Ω is the symplectic structure on the space of solutions $\mathcal{S}^\mathbb{C}$. Define the Hilbert space, \mathcal{H} , representing physical solutions by selecting a subspace of solutions such that: (i) The ‘‘inner product’’ (scare quotes because it is not positive definite on $\mathcal{S}^\mathbb{C}$) $(y_1, y_2) = -i\Omega(\bar{y}_1, y_2)$ is positive definite on \mathcal{H} , (ii) $\text{span}(\mathcal{H}, \bar{\mathcal{H}}) = \mathcal{S}^\mathbb{C}$, and (iii) for all $z_1 \in \mathcal{H}$ and $z_2 \in \bar{\mathcal{H}}$, $(z_1, z_2) = 0$. Importantly, there will many choices of \mathcal{H} that satisfy these conditions. The Hilbert space of the full QFT will then be $\mathcal{F}_S(\mathcal{H})$, the symmetrised Fock space constructed from \mathcal{H} .

¹³Given a positive frequency subspace with respect to a Killing vector field, time translating any state along the Killing vector field will recover a positive frequency state. Thus, the energy of the particle will always be positive when transformed by a time translation symmetry, as we desire for a particle interpretation. See Halvorson and Clifton (2002, pp. 3-4) for a brief discussion.

¹⁴See Arageorgis, Earman, and Ruetsche (2002) for a challenge to the possibility of formulating unitarily implementable dynamics for quantum field theories on generic, curved spacetimes.

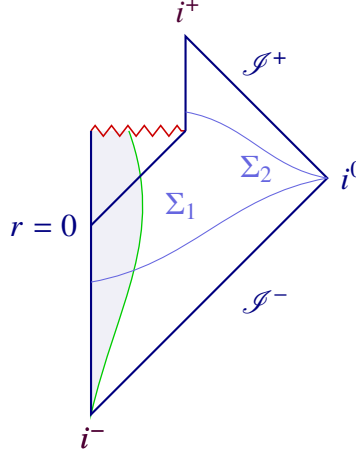


Figure 2: The conformal diagram representing the evaporation of a Schwarzschild black hole formed by collapse. The mass of the black hole is shrinking over time, and after the evaporation, the spacetime is locally Minkowski. Neither Σ_1 nor Σ_2 is a Cauchy surface.

without a full theory of quantum gravity, this interaction between the spacetime and the radiation can't be accounted for in full rigour. One can approximate the backreaction in two ways: either by modelling the radiation as a flux going out to infinity and using a conservation law to infer a flux down over the horizon, or using the semi-classical Einstein equation, $G_{ab} = 8\pi\langle T_{ab} \rangle$. The consensus view is that Hawking radiation implies a black hole loses mass-energy on pain of a “drastic violation of energy conservation” (Fredenhagen & Haag, 1990, pp. 282).¹⁵

The above approaches imply the black hole will radiate away all of its mass in finite time. The semi-classical approximation is expected to break down at late times, when the radius of the black hole is of the order of the Planck length. Beyond this point there much disagreement about the description of Hawking radiation and evaporation. However, the consensus has varied very little from Hawking's original heuristic: “there is not much it can do except disappear altogether.” (Hawking, 1975, pp. 219) Thus, black holes are usually supposed to evaporate entirely, with the spacetime in the region after evaporation isometric to a region of Minkowski spacetime. The conformal diagram for this spacetime, which I call evaporation-Schwarzschild, is depicted in figure 2.

Since this spacetime will be central to my discussion, I will highlight a few important features of it. It is very different from collapse-Schwarzschild: the metric in the region exterior to the collapsing matter is not Schwarzschild, it is not globally hyperbolic, it does not admit a timelike Killing vector field and it has a naked singularity, among other things. A consequence of Hawking radiation is that collapse-Schwarzschild is the wrong spacetime to describe the target black hole; it is an idealization. Given Hawking radiation, an uncharged, non-rotating black hole should be described by evaporation-

¹⁵See Wald (1995, sec. 7.3) for a treatment of the energy flux approach, and Wallace (2018) for a general overview of results in the semi-classical Einstein equation approach.

Schwarzschild. And yet, Hawking derived the eponymous radiation in collapse-Schwarzschild, in spite of the fact that many properties of collapse-Schwarzschild that are used in Hawking's derivation do not hold in evaporation-Schwarzschild; this threatens an essential idealization. Thus we arrive at the paradox discussed in the introduction. The rest of this paper defends the claim that throwing away the ladder of collapse-Schwarzschild really leads to inconsistency.

3 Idealization Paradox in Hawking's Derivation

3.1 Sketch of Hawking's Derivation

To understand exactly what goes wrong for Hawking's derivation in evaporation-Schwarzschild, we will need a more precise account of it. I give this here, stripped of unnecessary details.

In outline, we wish to compare the modes of a quantum field in the distant past with those in the distant future. Consider collapse-Schwarzschild spacetime¹⁶ containing a massless complex-valued scalar quantum field Φ (obtained as discussed above). Let $\{f_i\}$ be a complete basis of solutions, so that we may write: $\Phi = \sum_i \{f_i \mathbf{a}_i + \bar{f}_i \mathbf{a}_i^\dagger\}$, where \mathbf{a}_i and \mathbf{a}_i^\dagger are the annihilation and creation operators corresponding to the i th solution. We choose $\{f_i\}$ to be positive frequency solutions with respect to a time parameter defined by a timelike Killing vector field asymptotically close \mathcal{I}^- .

We can also describe Φ as a decomposition into solutions at \mathcal{I}^+ and on the event horizon \mathcal{H}^E . At \mathcal{I}^+ we can again form a Hilbert space generated by positive frequency solutions, $\{p_i\}$, with respect to a time parameter defined by a timelike Killing vector field on \mathcal{I}^+ . The modes on \mathcal{H}^E play no role in the derivation. \mathbf{b}_i , \mathbf{b}_i^\dagger are the annihilation and creation operators for the p_i modes. Because the spacetime is globally hyperbolic, we can express $\{p_i\}$ and \mathbf{b}_i as linear combinations of $\{f_i\}$ and $\{\bar{f}_i\}$ and \mathbf{a}_i and \mathbf{a}_i^\dagger respectively,

$$p_i = \sum_j \{\alpha_{ij} f_j + \beta_{ij} \bar{f}_j\}, \quad \mathbf{b}_i = \sum_j \{\bar{\alpha}_{ij} \mathbf{a}_j - \bar{\beta}_{ij} \mathbf{a}_j^\dagger\} \quad (1)$$

Stipulate that the field is in the state $|0_-\rangle$ defined as the vacuum state at early times: $\mathbf{a}_i |0_-\rangle = \mathbf{0}$ for all i . On \mathcal{I}^+ , $\mathbf{b}_i^\dagger \mathbf{b}_i$ has expectation value

$$\langle 0_- | \mathbf{b}_i^\dagger \mathbf{b}_i | 0_- \rangle = \sum_j |\beta_{ij}|^2 \quad (2)$$

which will be non-zero because we have different Killing vector fields defining our Hilbert spaces. Thus, to determine the expected number of particles in each mode, one needs to determine the coefficients β_{ij} .

For this calculation, Hawking writes the modes of Φ in terms of advanced and retarded

¹⁶Hawking also derived the radiation for charged, rotating black holes, but I focus on the simplest case.

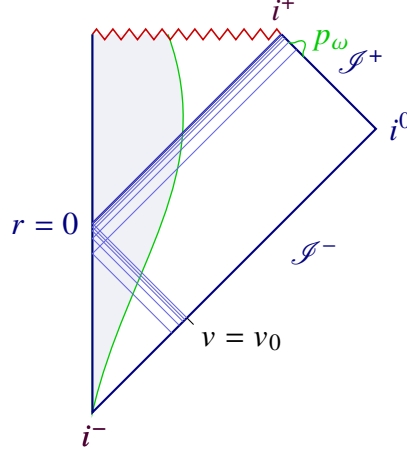


Figure 3: Conformal diagram used to visualise the mapping of modes of a quantum field on \mathcal{I}^+ to modes on \mathcal{I}^- , which pass through the non-stationary region of collapsing matter.

Eddington-Finkelstein coordinates:

$$v = t + r_*, \quad u = t - r_*, \quad r_* = r + 2m \log \left| \frac{r}{2m} - 1 \right| \quad (3)$$

Hawking considers a mode p_i on \mathcal{I}^+ at late retarded time u of frequency ω , defined with respect to retarded time, $p_\omega(u)$. He propagates this mode back along the event horizon through the non-stationary region of the collapsing star onto the \mathcal{I}^- (see figure 3). The form of the mode on \mathcal{I}^- is determined by connecting the mode to the event horizon by a null vector normal to the horizon, and parallel transporting this vector onto \mathcal{I}^- .¹⁷ From the form of the mode on \mathcal{I}^- , one can read off the β coefficients. Thus we arrive at Hawking's discovery: the expected particle number at frequency ω at \mathcal{I}^+ is that of a black body with temperature, in geometric units, of $\frac{\kappa}{2\pi}$, where κ is the surface gravity of the black hole. The black hole is seemingly radiating at what is now called the Hawking temperature.

Our task now is to investigate why this derivation cannot be carried out in evaporation spacetime. I begin by identifying a globally hyperbolic sub-spacetime of evaporation Schwarzschild that might plausibly admit Hawking's derivation. I then show that structure used in Hawking's derivation is not present in evaporation spacetimes and so Hawking's derivation is evaporation-inconsistent.

3.2 Hawking's Derivation Fails in Evaporation-Schwarzschild

Our first task is to find the region of evaporation-Schwarzschild in which to attempt to recover Hawking radiation. In quantum field theory on curved spacetimes, global hyperbolicity is nearly always assumed because this guarantees an initial value problem in the following sense: given an initial data

¹⁷In fact Hawking conducts the calculation on the past horizon of maximally extended Schwarzschild and argues that the conclusions would be the same on \mathcal{I}^- .

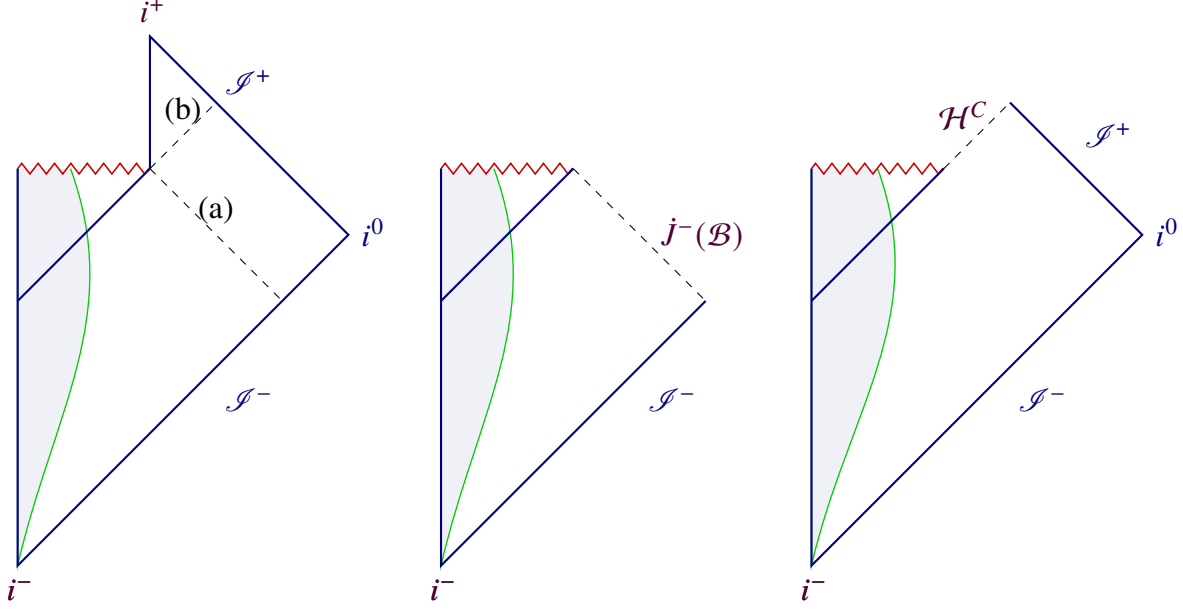


Figure 4: Left: The two globally hyperbolic regions of evaporation-Schwarzschild. The region below (a) is the causal past of the black hole, and the region below (b) is the MGHD of \mathcal{I}^- . Centre: The causal past of the black hole, $J^-(\mathcal{B})$. Right: The MGHD of \mathcal{I}^- , $D(\mathcal{I}^-)$.

surface in GR, there exists a unique (up to isometry) spacetime that is the maximal globally hyperbolic development (MGHD) of the data surface. The initial data surface will be a Cauchy surface for this spacetime, and determines the entire spacetime. Moreover, it is clear that derivations of Hawking radiation which map modes in the past to modes in the future (as Hawking's and Fredenhagen and Haag's do) will require global hyperbolicity. This is because the state of the field in the past must determine the state of the field in the future. However, evaporation-Schwarzschild is not globally hyperbolic. Therefore, we must find a region of evaporation-Schwarzschild which is globally hyperbolic and has sufficient spacetime structure to admit Hawking's derivation of Hawking radiation. There are two reasonable sub-spacetime regions of evaporation-Schwarzschild which are globally hyperbolic: a) the causal past of the black hole region ($J^-(\mathcal{B})$), or b) the MGHD of \mathcal{I}^- , ($D(\mathcal{I}^-)$). These two embedded regions are demarcated in figure 4.

Which is more suited to deriving Hawking radiation? It is MGHD \mathcal{I}^- . To see this, consider the causal past of the black hole spacetime. It is a spacetime such that, if a light ray were admitted at a point, it could reach the black hole before it evaporates completely. Near the evaporation event this is a tiny area, so we have deleted most of the spacetime we need for the derivation.

More technically, 'future null infinity' in the causal past of the black hole will be the boundary of the causal past, $\dot{J}^-(\mathcal{B})$. There will not be a timelike Killing vector field on this boundary; therefore, there will be no preferred time parameter with respect to which we can define a particle interpretation. This is because the boundary bisects the non-stationary exterior region of the spacetime. This also

means the boundary won't be asymptotically flat; instead, it ends at the naked singularity and so it will contain a region of arbitrarily large curvature. Clearly, the causal past of the black hole region is useless for deriving Hawking radiation.

MGHD \mathcal{I}^- on the other hand does not suffer these problems, and includes the portion of \mathcal{I}^+ where all Hawking radiation will propagate to. Therefore this is the appropriate globally hyperbolic spacetime region to use.¹⁸ So the question of this section is precisely stated as: which of the necessary assumptions for Hawking's derivation of Hawking radiation cannot be carried over into MGHD \mathcal{I}^- ?

There are two important differences between MGHD \mathcal{I}^- and collapse-Schwarzschild: the exterior solution is not Schwarzschild and the spacetime is not stationary. How do these changes affect the derivation? Firstly, Hawking's derivation makes use of ingoing and outgoing Eddington-Finkelstein coordinates, defined in equation (3), which are specified for a particular mass m . This constant mass term is unavailable in evaporation spacetimes. Instead, one must analyse how modes defined with respect to coordinates that cover MGHD \mathcal{I}^- behave on an evaporation metric, but nothing like this is carried out for Hawking's derivation.

Next, to calculate the form of modes on \mathcal{I}^- , Hawking exploits an isometry with the maximally extended Schwarzschild solution, and analyses modes that propagating onto the past horizon (see footnote 17). When the exterior solution is no longer Schwarzschild this isometry can not be used. Furthermore, and perhaps most strikingly, the failure of stationarity implies that the propagation of the modes back along the horizon will induce an evolution of the modes different to that calculated in collapse-Schwarzschild. Indeed, the normal null vector on the horizon which is used to compute the backwards evolution of the modes will have a different form in MGHD \mathcal{I}^- as compared with collapse-Schwarzschild, precisely because the metric is different and the horizon area is changing. Finally, the non-stationarity will affect the scattering of the modes by the gravitational field.

Admittedly, the model of evaporation used here, evaporation-Schwarzschild, is heuristic only and not generally believed to be a realistic model of black hole evaporation. One may wonder whether in more realistic models of black hole evaporation the problems listed here go away. It is in fact the opposite, things are worse in realistic models. For example, in explicitly computed models Schindler, Aguirre, and Kuttner (2020) show that, as well as the above worries still holding true, there is also no event horizon or Killing horizon for an evaporating black hole. Thus there will be no null vector normal to the horizon at all; the very structure Hawking uses to compute the form of the modes on \mathcal{I}^- is non-existent in evaporation spacetimes. So, in realistic evaporation models, more of the spacetime structure exploited by Hawking to derive the radiation is lost.

Hawking himself notes that the “negative energy flux will cause the area of the event horizon to decrease and so the black hole will not, in fact, be in a stationary state” (1975, p. 219). He accepts

¹⁸It can be shown that neither the causal past of the black hole nor MGHD \mathcal{I}^- is conformally equivalent to collapse-Schwarzschild, so proofs of the conformal equivalence of the Hawking temperature (e.g. Jacobson and Kang (1993)) do not help resolve the paradox.

this is a problem, but claims to have a solution, as one can approximate the black hole as “quasi stationary”. In section 6.2 I show that this approximation regime does not, in fact, save Hawking’s derivation because one cannot use the regime to recover the necessary global structure. Therefore, the problems remain.

A reader familiar with the vast literature of derivations of Hawking radiation may at this point be thinking of their preferred derivations, and be under the impression that they do not fall victims to the above challenges. I have no objection to such claims. Indeed I will present certain derivations as the best candidates currently available to resolve the paradox for Hawking’s derivation in section 6.3. Nevertheless, this is not a problem for my thesis as I am focused on particular derivations of Hawking radiation, in this case Hawking’s original derivation. Thus, given the amount of structure exploited by Hawking which does not carry over to MGHD \mathcal{I}^- , one must accept the conclusion that the derivation of Hawking radiation found in Hawking (1975) falls victim to the idealization paradox. That is to say, remarkably, Hawking’s derivation is evaporation-inconsistent!

Fredenhagen and Haag (1990) construct their derivation to avoid a different problematic assumption in Hawking’s derivation, the geometric optics approximation. Thus, Fredenhagen and Haag’s derivation is what most consider to be the watertight derivation. It is to this that I turn next.

4 Idealization Paradox in Fredenhagen and Haag’s derivation

Fredenhagen and Haag’s derivation is similar to Hawking’s in that it defines the state outside the black hole at some early time and maps this state to some state at late time. However, it differs in a few important respects. Firstly, the entire calculation is performed on the region of spacetime after the stellar matter has passed the event horizon. Secondly, they use a ‘detector’ at asymptotically late times to model the radiation. Thirdly, they perform the calculation by propagating the detector along the timelike Killing vector field in the exterior region. I sketch this derivation next, and in section 4.2 show that it is also evaporation-inconsistent.

4.1 Sketch of Fredenhagen and Haag’s Derivation

This derivation, like Hawking’s, takes place on collapse-Schwarzschild. The region exterior to the event horizon in Schwarzschild can be covered by the coordinates (t, r, θ, ϕ) , where we call t Schwarzschild-time, and define τ -time coordinates, (τ, r, θ, ϕ) where $\tau = t + r^* - r = v - r$, for v and r^* defined in (3). τ is approximately Schwarzschild-time near spacelike infinity, and becomes infinitely negative near the horizon. Let $r = r_s(\tau)$ define the surface of the collapsing star, with $r_s(0) = r_0$ the Schwarzschild radius, such that the star crosses the Schwarzschild radius at $\tau = 0$.

As before, let Φ be a massless complex-valued scalar quantum field which satisfies the covariant wave equation. Fredenhagen and Haag model a detector in a spacetime region O with an observable,

Q^*Q , which is the counting rate given by $\langle Q^*Q \rangle$, where $Q = \int \Phi(x)h(x)\sqrt{|g|}d^4x$ for a test function $h(x)$ that has support in O . They ‘place’ the detector at a large radius at the τ -time for which the collapsing star crosses the horizon (i.e. h has support around $(0, R, \theta_0, \phi_0)$ for $R \gg r_0$). The detector is then translated along the timelike Killing vector field of the Schwarzschild metric. We are interested in the counting rate of the detector at asymptotically late times (given by $Q_T(T \rightarrow \infty)$), as displayed in figure 5, in which a collapsing star and the time-translated detector are displayed in τ -time coordinates.

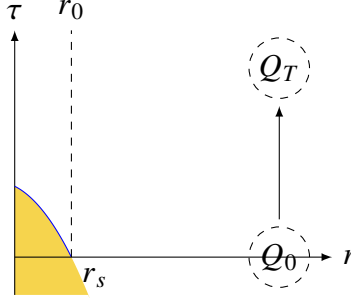


Figure 5: Set up of Fredenhagen and Haag’s derivation, in which a detector Q is propagated along the timelike Killing vector field to asymptotically late times, $Q_T(T \rightarrow \infty)$.

The counting rate is determined by the data on a Cauchy surface in the past of the late time detector. In the asymptotic limit, $T \rightarrow \infty$, the contributing data on the Cauchy surface decomposes into a sum of two wave packets, shifting asymptotically to $r \rightarrow \infty$ and $r \rightarrow r_0$. Wald (1995, pp. 159-162) explains this fact by noting that in maximally extended Schwarzschild, any mode in the region exterior to the black hole will decompose into modes on the past horizon and \mathcal{I}^- . Propagating this decomposition along Killing vector fields infinitely far will place modes infinitely close to the future horizon and spatial infinity, as depicted in figure 6. By isometry, we can draw the same conclusion for the exterior of collapse-Schwarzschild.¹⁹ This fact can also be seen as a consequence of the Schwarzschild potential pushing modes onto the horizon and out to infinity.

Assuming the state in the distant past is vacuum, the contribution to the counting rate from spatial infinity is zero. The contribution from the wave-packet that accumulates at the horizon is determined by the short-distance behaviour of the quantum field. Without concerning ourselves with the details, the authors assume the leading singularity in the short distance behaviour has a particular form that is implied by the Hadamard condition. Armed with this assumption, Fredenhagen and Haag show that the modes on the horizon contribute a thermal spectrum to the counting rate of the detector at asymptotically late times, with temperature given by the Hawking temperature.

¹⁹This inference is not in fact secure because the MGHD of the spacetime region exterior to the collapsing matter is not maximally extended Schwarzschild, but I ignore this difficulty here as it does not undermine Fredenhagen and Haag’s calculation but only Wald’s explanation of the behaviour of the wave-packet decomposition.

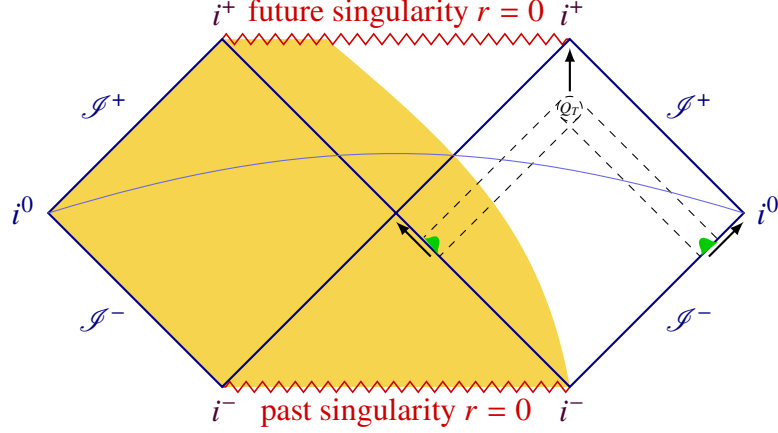


Figure 6: Decomposition of modes contributing to detector response of maximally extended Schwarzschild.

4.2 Fredenhagen and Haag’s Derivation Fails in Evaporation-Schwarzschild

We are again interested in whether we can deidealize this derivation: can the assumptions necessary to carry out Fredenhagen and Haag’s derivation be carried over to MGHD \mathcal{I}^- ?

The most prominent difficulty for the derivation is that it relies on the stationarity of the exterior metric. The detector is time translated along the Killing vector field of the Schwarzschild metric. This sends the modes on the Cauchy surface that contribute to the counting rate to spatial infinity and onto the horizon. Moreover, the behaviour of the mode decomposition as the detector is time translated is analysed on the maximally extended Schwarzschild spacetime and, following Wald, arises due to the global Killing field that is timelike in the exterior region.

In MGHD \mathcal{I}^- there are no such timelike Killing vector fields and there is no isometry with maximally extended Schwarzschild because the size of the black hole is changing. Indeed, there do not even exist approximate Killing vector fields on the entire spacetime, whatever notion of ‘approximate’ one might try to use. MGHD \mathcal{I}^- contains a large mass black hole at $\tau = 0$, and by the evaporation event it contains a negligible mass black hole. This is clearly a radical change and so the spacetime is in no sense stationary. The behaviour of modes under the time translation symmetry of collapse-Schwarzschild was the core of the derivation, and this is simply not available in evaporation-Schwarzschild.

In addition, there is a further difference between the two spacetimes relevant to Fredenhagen and Haag. The global time function on MGHD \mathcal{I}^- does not extend to infinity into the future, whereas it is future infinite on collapse-Schwarzschild. The lack of a future-infinite time coordinate is a problem for Fredenhagen and Haag’s derivation because, whereas Hawking’s asymptotic time assumption was realised by future null infinity, Fredenhagen and Haag translate their detector along a timelike worldline. Every timelike worldline will reach the Cauchy horizon of MGHD \mathcal{I}^- in a finite parameter

distance, so one cannot take the asymptotic time limit. This limit was essential to Fredenhagen and Haag's derivation as it pushed the modes asymptotically close to the horizon, forcing them into the trans-Planckian regime. Fredenhagen and Haag can then describe the modes by their short distance behaviour. Without this limit, we cannot be sure of the derivation.

The Fredenhagen and Haag derivation cannot be carried out, in any obvious fashion, in evaporation-Schwarzschild. The idealization paradox thus applies to this approach as well: it too is evaporation-inconsistent. I now turn to algebraic approaches.

5 Idealization Paradox in Algebraic Approaches

5.1 Sketch of the Algebraic Derivation

I will not go into any sort of detail in the sketch of algebraic approaches, as they are on the one hand very mathematically heavy, but on the other very conceptually simple. Algebraic approaches function by showing a particular state is the uniquely natural stationary Hadamard vacuum state on the collapse-Schwarzschild spacetime, and that this state is thermal at the Hawking temperature at future null infinity.

Algebraic QFT begins with a $*$ -algebra of observables \mathcal{A} . A state ω is a completely positive map from \mathcal{A} to \mathbb{C} , $\omega : \mathcal{A} \rightarrow \mathbb{C}$. For self-adjoint operators the map is real valued. We fix states by demanding they obey certain conditions, such as being vacuum.²⁰ Conversely, we can discover facts about states by assessing what conditions they obey, for example a state is thermal with respect to a given Hamiltonian if it obeys the KMS condition.²¹ As usual, we demand that physical states are Hadamard. Finally, one can define the algebra of observables for a scalar field by demanding that the functions used to smear the observables solve the covariant-wave equation.

One finds (Dimock and Kay (1987), Dappiaggi et al. (2011)) that the uniquely natural stationary Hadamard vacuum state on the collapse-Schwarzschild spacetime is the Unruh vacuum. The Unruh vacuum has the property of having no particles near \mathcal{I}^- , but being thermal at the Hawking temperature near \mathcal{I}^+ , with a flux going out to infinity. Thus, one claims that the black hole is emitting Hawking radiation.

This sketch is sufficient to analyse the idealization paradox for algebraic approaches, to which I turn now.

²⁰In curved spacetimes, the algebraic approach defines a vacuum as a state that is Gaussian and pure (see Kay and Wald (1991) for details).

²¹See Bratelli and Robinson (1982, p. 13) for a definition.

5.2 Algebraic Approaches Fail in Evaporation-Schwarzschild

Algebraic approaches are the most mathematically rigorous formulation of Hawking radiation. However, they clearly fail to survive the move to evaporation-Schwarzschild, or MGHD \mathcal{S}^- .

The spacetime we are now interested in is not collapse-Schwarzschild, and not even approximately collapse-Schwarzschild. Therefore, the proof of the unique naturalness of the Unruh vacuum simply does not apply; the Unruh vacuum is uniquely natural on collapse-Schwarzschild, with no implication for the uniquely natural vacuum state on MGHD \mathcal{S}^- . Moreover, given that one condition on the Unruh vacuum is that it is stationary, and collapse-Schwarzschild is not stationary, clearly the Unruh vacuum will be the inappropriate vacuum state for MGHD \mathcal{S}^- . We can thus conclude that the idealization paradox applies to the algebraic approaches.

To conclude, we have three different derivations, each of increasing mathematical rigour, and each with open questions about how they can actually claim to be establishing Hawking radiation in physically realistic models.

6 Paths Toward a Resolution

Physics uses idealizations all the time. The idealization used in the derivations here is only particularly striking because it leads to a paradox, rendering the argument each derivation presents for Hawking radiation inconsistent. This paradox is clearly unacceptable, and so we should find a resolution. One aim of this paper, the aim taken up in this section, is to categorise and assess solutions to the idealization paradox. The most natural resolutions are those that deidealize the derivations, to show how they can proceed in evaporation spacetimes. Three sub-categories of deidealization solutions are presented below:

- Quantum Gravity (section 6.1)
- Approximation Regime (section 6.2)
- Essential Structure (section 6.3)

The first argues that quantum gravity is needed to describe black hole evaporation and thus resolve the paradox. The second looks to find an approximation regime between collapse-Schwarzschild and evaporation-Schwarzschild. Specifically, I formalise and analyse an approximation regime suggested in Hawking (1975). The third argues that one can weaken the assumptions of the derivations, such that each derivation can derive Hawking radiation whilst assuming only some essential spacetime structure that is present in both evaporation and non-evaporation spacetimes.

I find that quantum gravity holds no prospects for resolving the paradox. I find Hawking's approximation regime achieves varying degrees of success for the different derivations, but even where there are hints of success more work is needed. Finally, I find that essential structure derivations constitute a very fruitful research direction which has already been taken up in Visser (2003); Barcelo,

Liberati, Sonogo, and Visser (2011b, 2011a). Indeed, this work already points towards deep lessons about the nature of Hawking radiation.

The derivations analysed in this paper, as discussed in the introduction, are not exhaustive. So, plausibly, other derivations don't face the paradox (indeed I will discuss some examples in section 6.3). Given a paradox-free derivation, the consistency conjecture will be true, and so the phenomenon of Hawking radiation will be insulated from the paradox.

However, this still leaves us with an idealization paradox for at least some derivations, and this in turn leaves an important and unanswered question: why is such a scientifically revolutionary piece of physics, Hawking (1975), inconsistent? How did it lead to Hawking radiation, when we can't mesh it with physically realistic black holes? Although, according to many, the context of discovery need not be objective and rational (c.f. Popper (1959, sec. 2)), it still seems highly unlikely that Hawking's derivation 'got lucky' and so an explanation of its success should be sought.

I do not consider here resolutions which may be collected under the name deidealization pessimism, examples of such views include: embracing evaporation-inconsistent derivations as essential idealizations (aligning with Batterman (2002, 2005, 2011)), and denying the phenomenon of either black hole evaporation or Hawking radiation. Such approaches would resolve the paradox, but offer a somewhat pyrrhic victory by respectively rejecting either Earman's principle, or the consensus in black hole physics.²² Instead, the categories I propose below (in sections 6.1, 6.2 and 6.3) help to distinguish different ways a derivation may be deidealized to avoid the paradox.

6.1 Quantum Gravity

It is widely believed that a quantum theory of gravity will resolve the black hole information paradox.²³ This is because the consensus in the physics community is that a quantum theory of gravity will be required to describe the final stages of black hole evaporation (e.g. Rovelli and Vidotto (2014)). Moreover, it is often claimed that the early stages of black hole evaporation also cannot be fully described without a quantum theory of gravity, as we can't accurately describe the backreaction of Hawking radiation on the metric. Therefore, a reasonable first suggestion is to expect quantum gravity to resolve the idealization paradox. However, I argue this proposal cannot succeed.

The central idea of a quantum gravity resolution to the idealization paradox is that the physics of spacetimes and Hawking radiation occurs in the semi-classical limit, whereas black hole evaporation occurs in a full quantum gravity description. One would argue that this allows one to reject the Inconsistency Claim, premise 3 of the idealization paradox presented in the introduction, which asserts that black hole evaporation leads to the rejection of assumptions required for the derivation of Hawking radiation. In order to reject the Inconsistency Claim, one may argue that because black hole

²²Wallace (2018, 2019) reviews the arguments in favour of this consensus.

²³For taxonomies of such proposals see Belot et al. (1999); Unruh and Wald (2017).

evaporation is a quantum gravity phenomenon, it is not describable in the semi-classical limit and as such tells us nothing about the properties of spacetime in the semi-classical limit. Thus one cannot infer from evaporation to the breakdown of the semi-classical limit spacetime properties required for Hawking radiation. Thus, by acknowledging the need for a quantum theory of gravity to describe black hole evaporation, we can escape the paradox.

Unfortunately, quantum gravity does not license us to reject the Inconsistency Claim. To see this, note that any quantum gravity theory of black hole evaporation must be able to represent: i) a black hole of given mass-energy, and ii) the mass-energy of a black hole being reduced in the process of evaporation. If the mass-energy of a black hole is not reducing then one cannot claim to be describing black hole evaporation, it is some other phenomena. This is certainly a possibility, but such a theory would constitute evaporation scepticism by claiming Hawking radiation does not lead black holes to lose mass-energy.

Given these minimal representational requirements, the state in our quantum theory of gravity will represent a black hole of mass m_1 in the semi-classical limit at some earlier time, and a black hole of mass m_2 in the semi-classical limit at some later time, where $m_1 > m_2$. This immediately violates stationarity, one of the properties used in the derivations of Hawking radiation discussed here. Therefore, even a completely quantum gravity model of evaporation implies the breakdown of properties required for derivations of Hawking radiation in the semi-classical limit. Hence we are not licensed to reject the Inconsistency Claim.

Why does the black hole information paradox admit a quantum gravity resolution whereas the idealization paradox does not? The difference is that there is no black hole information paradox until the evaporation event²⁴ because only then does one have to accept the information has vanished from the universe. Moreover, there is no minimal representational requirement on the evaporation event so we cannot anticipate any aspect of the quantum gravity description. On the other hand, the idealization paradox arises without the need to consider the evaporation event, due to the failure of properties in the entire exterior region such as stationarity. We can then impose our minimal condition on evaporation far before the evaporation event, and this leads to the paradox.

6.2 Approximation Regime

Perhaps the derivations considered here can be carried out in some appropriate approximation regime: One would find some spacetime region in collapse-Schwarzschild which looks approximately like some corresponding region of evaporation-Schwarzschild. One could then hope to carry out the derivation using this approximating region of collapse-Schwarzschild, then infer the derived radiation back onto evaporation-Schwarzschild. Thus, one would derive the existence of Hawking radiation in the evaporation spacetime. Hawking (1975, p. 219) proposed such a resolution to the paradox: “it

²⁴In the traditional sense, though not in the Page-time paradox sense; see Wallace (2020).

is a reasonable approximation to describe the black hole by a sequence of stationary solutions and to calculate the rate of particle emission in each solution.”

The regime is justified as follows: The rate of change of the mass of the black hole will (for masses larger than the Planck mass) be much slower than the time taken for light to propagate to a region that can be modelled as approximately flat.²⁵ Thus, one can approximate the variable mass black hole spacetime as a sequence of stationary regions and calculate the rate of particle emission in each solution, avoiding the non-stationarity issues.

This is a very intuitive picture if one imagines a black hole as a compact three-dimensional object that evolves in time and produces Hawking radiation via a local mechanism. However, as we have seen, the derivations of Hawking radiation discussed above use global spacetime structure, including the propagation of modes through the collapsing matter region (Hawking, 1975) and timelike Killing vector fields with an future infinite time parameter (Fredenhagen & Haag, 1990). Consequently, we should not assume that the slow rate of evaporation is sufficient to guarantee the derivations are unaffected; indeed, to do so would be negligent of philosophers of physics seeking to understand the derivations, idealizations, and phenomena at hand.

To give Hawking a more charitable treatment, let me propose a more promising way to formalise this approximation. We want to identify regions of evaporation-Schwarzschild which are approximated by regions of collapse-Schwarzschild so that the derivations can be carried out using collapse-Schwarzschild. The best candidate regions are the parts of spacetime which Hawking radiation propagates through as it escapes to future null infinity. Figure 7 illustrates the structure of the regions (following Wald (1995, p. 178)). We model photons carrying energy away from the black hole to \mathcal{I}^+ and a negative energy flux propagating over the horizon. This is symbolised by two red arrows emerging from a single point, one pointing over the event horizon, the other out to future null infinity. Two such photon emission events are displayed in each conformal diagram in figure 7. The shaded region between the two photon emission events in evaporation-Schwarzschild is quasi-stationary. Thus the corresponding shaded region of the collapse-Schwarzschild spacetime of mass m is approximately isometric to the shaded region of evaporation-Schwarzschild, where m is the mass of the black hole according to the quasi-stationary region.²⁶

I denote such a quasi-stationary region of evaporation-Schwarzschild as \mathcal{R}_{QS} and a corresponding approximately isometric stationary region of a collapse-Schwarzschild spacetime as \mathcal{R}_S . The regime

²⁵Initial work on modelling ‘infinity’ at a finite distance so asymptotic flatness can be defined for realistic sub-systems of the Universe can be found in Ellis (2002, sec. 5). Such modelling frameworks will likely be helpful for rigorously analysing approximation regimes such as the one proposed here. However, finite-infinity models are tangential to our current concern because I accept the standard presumption of the literature that astrophysical black holes are well modelled by collapse-Schwarzschild. If one demands finite-infinity models, it is then necessary to show how derivations of Hawking radiation look in these deidealized models, and what the relationship of such finite-infinity models to the models discussed in this paper is.

²⁶The mass of the slices will have to be modelled by the Bondi mass, as this is defined at future null infinity whereas the ADM mass can only be defined at spacelike infinity.

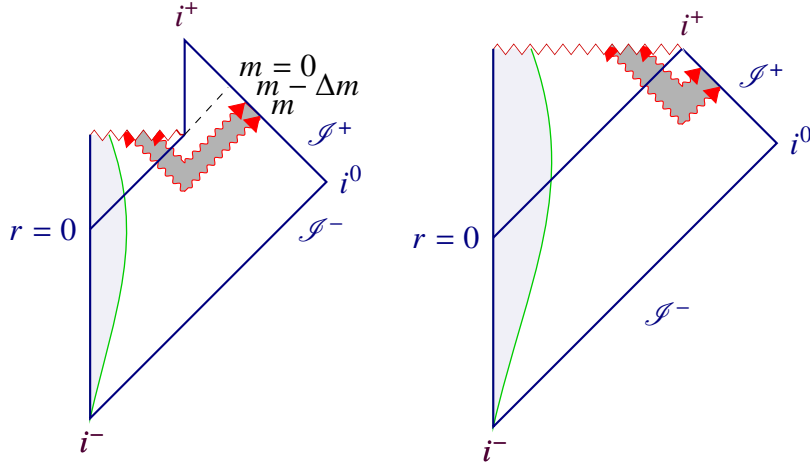


Figure 7: The two shaded regions are approximately isometric.

will work by using \mathcal{R}_S instead of \mathcal{R}_{QS} to derive Hawking radiation. One then infers the same result, to some degree of approximation, in the approximately isometric \mathcal{R}_{QS} . Repeating this for every \mathcal{R}_{QS} should describe the Hawking effect in evaporation-Schwarzschild.

This regime will face two central problems. First, for each of the derivations \mathcal{R}_S will have insufficient structure to derive Hawking radiation because it is a smaller extendable subspacetime. Thus it will be necessary to use the MGHD of \mathcal{R}_S . Problematically, although \mathcal{R}_{QS} and \mathcal{R}_S are approximately isometric, the corresponding MGHD for each will be very different. So although it is clear we can use \mathcal{R}_S to draw approximately correct conclusions about \mathcal{R}_{QS} , it is less clear that we can use the MGHD of \mathcal{R}_S to draw approximately correct conclusions about \mathcal{R}_{QS} . One must therefore justify using the MGHD of \mathcal{R}_S rather than only the approximately isometric region, \mathcal{R}_S , despite the different global structure.

Second, even if one can justify using the MGHD of \mathcal{R}_S to draw inferences about \mathcal{R}_{QS} , neither \mathcal{R}_{QS} nor \mathcal{R}_S contain Cauchy surfaces for evaporation-Schwarzschild or collapse-Schwarzschild respectively. This is obvious in figure 7 where, for example, a massive particle can travel from i^- and reach the singularity and never record data on the quasi-stationary surface. The same is true for \mathcal{R}_S in collapse-Schwarzschild. This means that neither \mathcal{R}_{QS} nor its approximately isometric stationary sibling \mathcal{R}_S determine the entirety of their respective spacetimes. In fact, the past domain of dependence for \mathcal{R}_{QS} does not extend outside of \mathcal{R}_{QS} , and so the past is significantly underdetermined.

Why can't we just select a region which does contain a Cauchy surface? Because this region would not be isometric, even approximately, to any region of collapse-Schwarzschild, and so we won't be able to use approximation to justify performing the derivation on collapse-Schwarzschild and transferring the result of the derivation back over to evaporation-Schwarzschild. Given this, let

us see how each of the derivations fail.

Consider Hawking’s derivation: it depends on global spacetime structure in the sense of an infinite past prior to collapse that is stationary, and an infinite future after collapse that is stationary, and a non-stationary intervening period. He writes: “To understand how the particle creation can arise from mixing of positive and negative frequencies, it is essential to consider not only the quasi-stationary final state of the black hole but also the time-dependent formation phase.” (Hawking, 1975, p. 204) \mathcal{R}_S contains none, or at most very little, of this requisite structure. For example, any given \mathcal{R}_{QS} need not intersect the collapse region; indeed, the majority not, as demonstrated in figure 7. Thus, the approximately isometric region \mathcal{R}_S will also not intersect the non-stationary collapsing matter and so will have insufficient structure to carry out Hawking’s derivation.

In order to recover the necessary structure, one needs to justify moving from \mathcal{R}_S to a spacetime with the global structure of collapse-Schwarzschild, perform the derivation on the global structure, and then make inferences from the global derivation back to the slice. Even if we assume that the first problem discussed above is solved and so such an inference is permissible, the inference still fails because, as per the second problem above, \mathcal{R}_S does not contain a Cauchy surface for collapse-Schwarzschild. Therefore, even if one could justify using the very different global structure to draw inferences about \mathcal{R}_{QS} , not enough of the global structure is included in the MGHD of \mathcal{R}_S to carry out Hawking’s derivation. Therefore, Hawking’s approximation regime fails for Hawking’s derivation.

Turning to Fredenhagen and Haag’s derivation: it was designed to not require the propagation of modes through the non-stationary collapse region, so the failure to recover this structure in \mathcal{R}_S is not problematic. However, the asymptotic time limit is not recovered in \mathcal{R}_S ; in fact the time over which the detector can be propagated is even shorter than in MGHD \mathcal{I}^- . Therefore, the modes cannot accumulate arbitrarily close to the horizon, as is needed in Fredenhagen and Haag’s derivation. However, perhaps the result is recovered approximately with this limited time evolution. Moreover, \mathcal{R}_S does determine the entire future of the spacetime, so if we can overcome the first problem above and justify using the MGHD of \mathcal{R}_S , we can in fact recover the asymptotic time-limit.

Hence, the approximation regime holds reasonable promise of succeeding for Fredenhagen and Haag’s derivation. Nonetheless, it needs to be shown that either the use of the MGHD of \mathcal{R}_S is justified, or \mathcal{R}_S admits sufficiently long stationary worldlines to allow modes to accumulate sufficiently close to the horizon that the singularity structure of the quantum field will dominate. Neither of these are trivial, but neither seems implausible either.

Finally, the algebraic approach: one begins by restricting the algebra of observables on collapse-Schwarzschild of mass m to an algebra on \mathcal{R}_S . One then infers the vacuum state on this algebra by restricting the Unruh vacuum to \mathcal{R}_S . One then claims that the algebra and vacuum state on \mathcal{R}_{QS} is approximately that of \mathcal{R}_S .

The central challenge for this approach is again the failure of \mathcal{R}_S to be Cauchy. This means

the uniqueness of the state on \mathcal{R}_S will probably not hold. Without uniqueness, we can't guarantee the state on \mathcal{R}_S is the restriction of the state on collapse-Schwarzschild. Moreover, the states on each \mathcal{R}_{QS} must be smoothly joined together, and therefore one needs to understand how the approximation changes the state, if only slightly.²⁷

I do not claim that I have exhausted the possibilities and difficulties for Hawking's proposal. Nor do I claim that Hawking's proposal exhausts the possible approximation regimes. I simply claim that, as of yet, this approach hasn't been completely worked out for any of the derivations. Moreover, if the approximation regime is worked out for one derivation, say Fredenhagen and Haag's, then the idealization paradox remains for the others, and so interesting open questions remain. The goal of this section has been to emphasise that the inference from that fact that evaporation is slow to the claim that the derivations go through approximately unaffected is non-trivial.

In the next section I consider whether we can weaken the premises of the derivations to deidealize them.

6.3 Essential Structure

The idealization paradox arises because some derivation uses a set of properties X with which to derive Hawking radiation, and then one finds that evaporation spacetimes don't instantiate the set of properties X . However, suppose that one could show that the derivation in fact did not require the complete set of properties X but only some subset of X , call it set Y , the essential structure. Suppose further that evaporation spacetimes could instantiate the essential structure Y . Then the inconsistency would be resolved. Moreover, the essential structure that goes into deriving Hawking radiation would have been identified, and the surplus structure is stripped away.

The task of identifying this essential structure is undertaken in Visser (2003); Barcelo et al. (2011a, 2011b). Visser (2003) argues that only three features are required for a derivation of Hawking radiation: an apparent horizon, non-zero surface gravity of the apparent horizon, and slow evolution. Therefore, using these as the set of properties Y could potentially resolve the idealization paradox. Going further, Barcelo et al. (2011a, 2011b) argue that Hawking-like radiation will occur whenever there is a continuous function mapping an affine parameter on future null infinity to that on past null infinity and the 'adiabatic condition' is satisfied.²⁸ In Barcelo et al. (2011a) the authors show how these conditions, with added assumptions about the QFT, can be used to derive the Bogoliubov coefficients, making explicit the relationship between their minimal conditions and Hawking's derivation of Hawking radiation.

Deidealization via essential structure derivations is strikingly different to that via the approximation regime. Whereas Hawking sought to find stationary structure within a non-stationary spacetime,

²⁷Work has begun to formulate algebraic QFT on non-globally hyperbolic spacetimes, e.g. Janssen (2022).

²⁸This is essentially a slow evolution condition, for details see Barcelo et al. (2011a).

these derivations do away with the need for quasi-stationary regions, and instead provide a derivation which would be successful on the global structure of an evaporation spacetime. Both Visser (2003) and Barcelo et al. (2011a, 2011b) require the black hole to evolve slowly, but they do not use this slow evolution to approximate stationarity. By using a different deidealization method, different lessons are drawn. For example, given the very minimal structure used in these derivations, it can be argued that they point towards a kinematic interpretation of Hawking radiation, *contra* the dynamical picture given in Hawking (1975). Moreover, these derivations don't require an event horizon or Killing horizon to form. Similar lessons to those from an approximation regime can be learned here also, for example the spectrum derived is only approximately thermal, and the spectrum can be derived away from the asymptotic future (i.e. before the retarded time coordinate goes to infinity).

I do not claim that these derivations face no difficulties, but only that they are very promising candidates for resolving the idealization paradox. A full analysis will be carried out in the future of the project. There is also a semantic issue of what one takes to be the referent of 'Hawking radiation' which I ignore here, emphasising only that resolutions to the paradox modify: i) what one takes to be required for something like Hawking radiation to occur, and ii) what is observed at \mathcal{I}^+ . On a cautious note, it is not clear that one can distinguish between radiation due to the Unruh effect and radiation due to the Hawking effect with these derivations. Although the Unruh effect and Hawking radiation are closely related phenomena, they are not the same (Earman, 2011). If one cannot distinguish between the two a derivation may have insufficient structure. However, this does not seem to me a serious obstacle to these derivations resolving the paradox, but rather an obstacle to the full interpretation of the Hawking effect.

The papers discussed here are not the only candidates for essential structure resolutions. Quantum tunnelling approaches, for example Parikh and Wilczek (2000), give a local dynamical account of Hawking radiation. A resolution to the paradox by these derivations would tell a different story. Firstly, they would retain a dynamical ontology for Hawking radiation. Secondly, they would point to lessons about the encoding of local structure by global structure in semi-classical gravity. Unpacking this second point, the definition of a black hole is global and Hawking's derivation is global, but if a resolution of the paradox along the lines of a local dynamical account is the correct one, we might learn that this global structure is a red herring, and it just encodes local structure that in ways that are, at times, opaque.

The goal of this section has been three-fold: 1) To highlight the differences between different deidealization strategies, 2) To emphasise there is an alternative to approximation regimes which make an inference from slow-evaporation to unaffected derivations, 3) To highlight the importance of research programmes such as that undertaken by Visser (2003); Barcelo et al. (2011a, 2011b). I do not claim that the paradox is definitely solved, or even necessarily solved by an essential structure deidealization, but rather that this is a promising option with many lessons to be learnt.

Summarising, deidealization can follow multiple different routes and these routes have varying degrees of success. Indeed, the success of a particular deidealization need not be homogeneous across derivations. I only take the quantum gravity route to be completely impotent. Hawking's approximation regime fails for Hawking's derivation, but prospects for success are better for Fredenhagen and Haag's derivation, and other approximation regimes may fare better. The essential structure research programme is very promising, in particular for deidealizing Hawking's derivation. The lessons we draw from these varying approaches to deidealization depend upon the type of deidealization, and the details of how the deidealization operates.

In the final section I consider what options are available if one thinks, against the presumed consensus, that one cannot deidealize the derivations presented here.

7 Conclusion

"Paradoxes are just the scar tissue. Time and space heal themselves up around them and people simply remember a version of events which makes as much sense as they require it to make." – Douglas Adams, *Dirk Gently's Holistic Detective Agency*

I have argued that Hawking's derivation of Hawking radiation, along with Fredenhagen and Haag's and the algebraic approach, are all evaporation-inconsistent. They are carried out on collapse-Schwarzschild but cannot be carried out on evaporation-Schwarzschild. By throwing away the space-time used to derive the phenomenon, we throw away the very ladder we are standing on, and come tumbling back to inconsistency. There are reasonable (and some unreasonable) paths towards deidealizing the derivations involved, and thus reason to believe that the paradox is just scar tissue from the messy process of scientific development. Presumably, there is a resolution along the lines of an approximation regime or essential structure derivation which will teach us why the inconsistent derivations worked so well, and what they really represent. If so, this will be another victory for the dispensabilists in the idealization literature, and a particularly striking one given that the idealizations used in Hawking radiation derivations were not simply false, but inconsistent. The differences between the different possible deidealizations emphasises the non-triviality of the deidealization project, and the variety of lessons that may be learnt. Most excitingly, deidealizing these derivations may remove the chaff from the conceptual framework of Hawking radiation and give a clear ontological picture of Hawking's eponymous discovery. Such lessons are the fruits of paying close attention to, and resolving, the idealization paradox; fruits won as reward for not settling with the unjustified inference from slow-evaporation to unaffected derivations. No matter what road we take, we are bound to learn something interesting.

Acknowledgements

I am very grateful to Bryan Roberts for his fantastic guidance and support in the development and refinement of this paper. I would also like to thank participants in the Philosophy of Physics Bootcamp for careful and helpful reading of an earlier draft, as well as simulating discussion on related topics. I am grateful to Erik Curiel whose generous reading and feedback improved this paper significantly. Finally, I am indebted to Bruno Arderucio, Jeremy Butterfield, Saakshi Dulani, Sam Fletcher, Henrique Gomes, Sean Gryb, Klaas Landsman, Miklós Rédei and Karim Thébault, and to audiences at: the University of Bristol Philosophy of Physics seminar, the Cosmology and Quantum Gravity Beyond Spacetime conference at Western University, the Golden Wedding of Black Holes and Thermodynamics conference, the Harvard Black Hole Initiative Foundations seminar and the annual conference of the Philosophy of Physics Group of the German Physical Society in Berlin.

References

- Arageorgis, A., Earman, J., & Ruetsche, L. (2002). Weyling the time away: the non-unitary implementability of quantum field dynamics on curved spacetime. *Studies in History and Philosophy of Science Part B: Studies in History and Philosophy of Modern Physics*, 33(2), 151–184.
- Barcelo, C., Liberati, S., Sonego, S., & Visser, M. (2011a). Hawking-like radiation from evolving black holes and compact horizonless objects. *Journal of High Energy Physics*, 2011(2), 1–30.
- Barcelo, C., Liberati, S., Sonego, S., & Visser, M. (2011b). Minimal conditions for the existence of a Hawking-like flux. *Physical Review D*, 83(4), 041501.
- Batterman, R. (2002). *The Devil in the Details: Asymptotic Reasoning in Explanation, Reduction, and Emergence*. New York: Oxford University Press.
- Batterman, R. (2005). Critical phenomena and breaking drops: Infinite idealizations in physics. *Studies in History and Philosophy of Science Part B: Studies in History and Philosophy of Modern Physics*, 36(2), 225–244. doi: <https://doi.org/10.1016/j.shpsb.2004.05.004>
- Batterman, R. (2011). Emergence, Singularities, and Symmetry Breaking. *Foundations of Physics*, 41(6), 1031–1050. doi: 10.1007/s10701-010-9493-4
- Batterman, R. (2017). Philosophical Implications of Kadanoff’s Work on the Renormalization Group. *Journal of Statistical Physics*, 167(3-4), 559–574. doi: 10.1007/s10955-016-1659-9
- Belot, G., Earman, J., & Ruetsche, L. (1999). The Hawking Information Loss Paradox: The Anatomy of Controversy. *The British Journal for the Philosophy of Science*, 50(2), 189–229. doi: 10.1093/bjps/50.2.189
- Bratelli, O., & Robinson, D. W. (1982). *Operator algebras and quantum statistical mechanics*.

- Butterfield, J. (2011). Less is different: Emergence and reduction reconciled. *Foundations of physics*, 41, 1065–1135.
- Curiel, E. (2019). The Many Definitions of a Black Hole. *Nature Astronomy*, 3, 27–34.
- Curiel, E. (2023). *The Hawking Effect, Its Desiderata and Its Discontents*. Retrieved from <https://www.youtube.com/watch?v=jFZ2HSkMTvY>
- Dappiaggi, C., Moretti, V., & Pinamonti, N. (2011). Rigorous construction and Hadamard property of the Unruh state in Schwarzschild spacetime. *Advances in Theoretical and Mathematical Physics*, 15(2), 355–447.
- Dimock, J., & Kay, B. S. (1987). Classical and quantum scattering theory for linear scalar fields on the Schwarzschild metric I. *Annals of Physics*, 175(2), 366–426. doi: [https://doi.org/10.1016/0003-4916\(87\)90214-4](https://doi.org/10.1016/0003-4916(87)90214-4)
- Dougherty, J., & Callender, C. (2016). *Black hole thermodynamics: More than an analogy?* Retrieved from <https://philsci-archive.pitt.edu/13195/>
- Earman, J. (2004). Curie's Principle and spontaneous symmetry breaking. *International Studies in the Philosophy of Science*, 18(2-3), 173–198. doi: 10.1080/0269859042000311299
- Earman, J. (2011). The Unruh Effect for Philosophers. *Studies in History and Philosophy of Science Part B: Studies in History and Philosophy of Modern Physics*, 42(2), 81–97. doi: 10.1016/j.shpsb.2011.04.001
- Ellis, G. F. R. (2002). Cosmology and local physics. *New Astronomy Reviews*, 46(11), 645–657. doi: 10.1016/S1387-6473(02)00234-8
- Fletcher, S. C. (2020). The principle of stability. *Philosophers' Imprint*, 20.
- Fletcher, S. C., Palacios, P., Ruetsche, L., & Shech, E. (2019). Infinite idealizations in science: an introduction. *Synthese*, 196, 1657–1669.
- Fredenhagen, K., & Haag, R. (1990). On the Derivation of Hawking Radiation Associated With the Formation of a Black Hole. *Commun. Math. Phys.*, 127, 273. doi: 10.1007/BF02096757
- Frigg, R. (2022). *Models and theories a philosophical inquiry*. Routledge.
- Frigg, R., & Hartmann, S. (2020). Models in Science. In E. N. Zalta (Ed.), *The Stanford encyclopedia of philosophy* (Spring 2020 ed.). Metaphysics Research Lab, Stanford University. <https://plato.stanford.edu/archives/spr2020/entries/models-science/>.
- Geroch, R. (1970). Domain of Dependence. *Journal of Mathematical Physics*, 11(2), 437–449. doi: 10.1063/1.1665157
- Gryb, S., Palacios, P., & Thébault, K. (2019). On the Universality of Hawking Radiation. *British Journal for the Philosophy of Science*, axz025. doi: 10.1093/bjps/axz025
- Halvorson, H., & Clifton, R. (2002). No place for particles in relativistic quantum theories? *Philosophy of Science*, 69(1), 1–28. doi: 10.1086/338939
- Hawking, S. W. (1975). Particle Creation by Black Holes. *Commun. Math. Phys.*, 43, 199–220.

- ([Erratum: *Commun. Math. Phys.* 46, 206 (1976)]) doi: 10.1007/BF02345020
- Hawking, S. W., & Ellis, G. F. (1973). *The large scale structure of space-time*. Cambridge university press.
- Jacobson, T., & Kang, G. (1993). Conformal invariance of black hole temperature. *Classical and Quantum Gravity*, 10(11), L201.
- Janssen, D. W. (2022). Quantum fields on semi-globally hyperbolic space-times. *Communications in Mathematical Physics*, 391(2), 669–705.
- Jones, N. J. (2006). *Ineliminable idealizations, phase transitions, and irreversibility* (Unpublished doctoral dissertation). The Ohio State University.
- Kay, B. S., & Wald, R. M. (1991). Theorems on the uniqueness and thermal properties of stationary, nonsingular, quasifree states on spacetimes with a bifurcate Killing horizon. *Physics Reports*, 207(2), 49-136. doi: [https://doi.org/10.1016/0370-1573\(91\)90015-E](https://doi.org/10.1016/0370-1573(91)90015-E)
- Landsman, K. (2013). Spontaneous symmetry breaking in quantum systems: Emergence or reduction? *Studies in History and Philosophy of Science Part B: Studies in History and Philosophy of Modern Physics*, 44(4), 379–394.
- Landsman, K. (2021). *Foundations of General Relativity: From Einstein to Black Holes*. Radboud University Press. doi: 10.54195/EFVF4478
- Manchak, J. B., & Weatherall, J. O. (2018). Paradox Regained? A Brief Comment on Maudlin on Black Hole Information Loss. *Foundations of Physics*, 48(6), 611–627. doi: 10.1007/s10701-018-0170-3
- Maudlin, T. (2017). *(Information) Paradox Lost*.
- Menon, T., & Callender, C. (2013). Ch-Ch-Changes Philosophical Questions Raised by Phase Transitions. In R. Batterman (Ed.), *The oxford handbook of philosophy of physics* (p. 189). OUP.
- Morrison, M. (2012). Emergent physics and micro-ontology. *Philosophy of Science*, 79(1), 141–166. doi: 10.1086/663240
- Norton, J. D. (2012). Approximation and Idealization: Why the Difference Matters. *Philosophy of Science*, 79(2), 207–232. doi: 10.1086/664746
- Page, D. N. (1994). Black hole information. In *Proceedings of the 5th canadian conference on general relativity and relativistic astrophysics* (Vol. 1, pp. 1–41).
- Palacios, P. (2019). Phase Transitions: A Challenge for Intertheoretic Reduction? *Philosophy of Science*, 86(4), 612–640. doi: 10.1086/704974
- Palacios, P. (2022). *Emergence and Reduction in Physics*. Cambridge University Press. doi: 10.1017/9781108901017
- Parikh, M. K., & Wilczek, F. (2000). Hawking Radiation As Tunneling. *Physical Review Letters*, 85(24), 5042–5045. doi: 10.1103/physrevlett.85.5042

- Popper, K. (1959). *The Logic of Scientific Discovery* (6th ed.). Routledge.
- Potochnik, A. (2017). *Idealization and the Aims of Science*. Chicago: University of Chicago Press.
- Prunkl, C. E. A., & Timpson, C. G. (2019). *Black Hole Entropy is Thermodynamic Entropy*.
- Raju, S. (2022). Lessons from the information paradox. *Physics Reports*, 943, 1-80. (Lessons from the information paradox) doi: <https://doi.org/10.1016/j.physrep.2021.10.001>
- Rovelli, C., & Vidotto, F. (2014). Planck stars. *International Journal of Modern Physics D*, 23(12), 1442026. doi: 10.1142/s0218271814420267
- Schindler, J. C., Aguirre, A., & Kuttner, A. (2020). Understanding black hole evaporation using explicitly computed Penrose diagrams. *Phys. Rev. D*, 101, 024010. doi: 10.1103/PhysRevD.101.024010
- Shech, E. (2018). Infinite Idealizations in Physics. *Philosophy Compass*, 13(9), e12514. doi: 10.1111/phc3.12514
- Shech, E. (2023). *Idealizations in Physics*. Cambridge, UK: Cambridge University Press.
- Unruh, W. G. (2014). Has Hawking radiation been measured? *Foundations of Physics*, 44, 532–545.
- Unruh, W. G., & Wald, R. M. (2017). Information loss. *Reports on Progress in Physics*, 80.
- Visser, M. (2003). Essential and inessential features of Hawking radiation. *International Journal of Modern Physics D*, 12(04), 649–661. doi: 10.1142/s0218271803003190
- Wald, R. M. (1975). On particle creation by black holes. *Communications in Mathematical Physics*, 45(1), 9–34.
- Wald, R. M. (1984). *General Relativity*. Chicago, USA: Chicago Univ. Pr. doi: 10.7208/chicago/9780226870373.001.0001
- Wald, R. M. (1995). *Quantum Field Theory in Curved Space-Time and Black Hole Thermodynamics*. Chicago, IL: University of Chicago Press.
- Wallace, D. (2018). *The case for black hole thermodynamics, Part I: phenomenological thermodynamics*.
- Wallace, D. (2019). The case for black hole thermodynamics part II: Statistical mechanics. *Studies in History and Philosophy of Science Part B: Studies in History and Philosophy of Modern Physics*, 66, 103-117. doi: <https://doi.org/10.1016/j.shpsb.2018.10.006>
- Wallace, D. (2020). Why Black Hole Information Loss Is Paradoxical. In N. Huggett, K. Matsubara, & C. Wüthrich (Eds.), *Beyond spacetime: The foundations of quantum gravity* (p. 209–236). Cambridge University Press. doi: 10.1017/9781108655705.013
- Wüthrich, C. (2019). Are black holes about information? In R. Dardashti, R. Dawid, & K. Thébault (Eds.), *Why trust a theory?: Epistemology of fundamental physics* (p. 202-224). Cambridge University Press. doi: 10.1017/9781108671224