

Wavenumber-explicit convergence of the hp -FEM for the full-space heterogeneous Helmholtz equation with smooth coefficients

D. Lafontaine*, E. A. Spence†, J. Wunsch‡

March 5, 2022

Abstract

A convergence theory for the hp -FEM applied to a variety of constant-coefficient Helmholtz problems was pioneered in the papers [35], [36], [15], [34]. This theory shows that, if the solution operator is bounded polynomially in the wavenumber k , then the Galerkin method is quasioptimal provided that $hk/p \leq C_1$ and $p \geq C_2 \log k$, where C_1 is sufficiently small, C_2 is sufficiently large, and both are independent of k, h , and p . The significance of this result is that if $hk/p = C_1$ and $p = C_2 \log k$, then quasioptimality is achieved with the total number of degrees of freedom proportional to k^d ; i.e., the hp -FEM does not suffer from the pollution effect.

This paper proves the analogous quasioptimality result for the heterogeneous (i.e. variable-coefficient) Helmholtz equation, posed in \mathbb{R}^d , $d = 2, 3$, with the Sommerfeld radiation condition at infinity, and C^∞ coefficients. We also prove a bound on the relative error of the Galerkin solution in the particular case of the plane-wave scattering problem. These are the first ever results on the wavenumber-explicit convergence of the hp -FEM for the Helmholtz equation with variable coefficients.

1 Introduction

1.1 Context

Over the last 10 years, a wavenumber-explicit convergence theory for the hp -FEM applied to the Helmholtz equation

$$\Delta u + k^2 u = -f \tag{1.1}$$

was established in the papers [35], [36], [15], [34]. This theory is based on decomposing solutions of the Helmholtz equation into two components:

- (i) an analytic component, satisfying bounds with the same k -dependence as those satisfied by the full Helmholtz solution, and
- (ii) a component with finite regularity, satisfying bounds with improved k -dependence compared to those satisfied by the full Helmholtz solution.

Such a decomposition was obtained for

- the Helmholtz equation (1.1) posed in \mathbb{R}^d , $d = 2, 3$, with compactly-supported f , and with the Sommerfeld radiation condition

$$\frac{\partial u}{\partial r}(x) - iku(x) = o\left(\frac{1}{r^{(d-1)/2}}\right) \tag{1.2}$$

as $r := |x| \rightarrow \infty$, uniformly in $\hat{x} := x/r$ [35, Lemma 3.5],

¹Department of Mathematical Sciences, University of Bath, Bath, BA2 7AY, UK, D.Lafontaine@bath.ac.uk

²Department of Mathematical Sciences, University of Bath, Bath, BA2 7AY, UK, E.A.Spence@bath.ac.uk

³Department of Mathematics, Northwestern University, 2033 Sheridan Road, Evanston IL 60208-2730, US, jwunsch@math.northwestern.edu

- the Helmholtz exterior Dirichlet problem where the obstacle has analytic boundary [36, Theorem 4.20],
- the Helmholtz interior impedance problem where the domain is either smooth ($d = 2, 3$) [36, Theorem 4.10], [34, Theorem 4.5], or polygonal [36, Theorem 4.10], [15, Theorem 3.2].

This decomposition was then used to prove quasioptimality of the hp -FEM applied to the standard Helmholtz variational formulation in [35], [36], [15], and applied to a discontinuous Galerkin formulation in [34]. Indeed, for the standard variational formulation (defined for the full-space problem in Definition 2.2 below) applied to the boundary value problems above, if the solution operator of the problem is bounded polynomially in k (see Definition 2.6 below), then there exist C_1, C_2 , and C_{qo} (independent of k, h , and p) such that if

$$\frac{hk}{p} \leq C_1 \quad \text{and} \quad p \geq C_2 \log k \quad (1.3)$$

then the Galerkin solution u_N exists, is unique, and satisfies

$$\|u - u_N\|_{H_k^1} \leq C_{\text{qo}} \min_{v_N \in V_N} \|u - v_N\|_{H_k^1},$$

where V_N is the hp approximation space and the norm $\|\cdot\|_{H_k^1}$ is the standard weighted H^1 norm (defined by (2.7) below). Since the total number of degrees of freedom of the approximation space is proportional to $(p/h)^d$, the significance of this result is that it shows there is a choice of h and p such that the Galerkin solution is quasioptimal, with quasioptimality constant (i.e. C_{qo}) independent of k , and with the total number of degrees of freedom proportional to k^d ; thus, with these choices of k and p , the hp -FEM does not suffer from the pollution effect [2].

Over the last few years, there has been increasing interest in the numerical analysis of the heterogeneous Helmholtz equation, i.e. the Helmholtz equation with variable coefficients

$$\nabla \cdot (\mathbf{A} \nabla u) + k^2 n u = -f; \quad (1.4)$$

see, e.g., [8], [3], [10], [18], [38], [21], [16], [29], [19]. However there do not yet exist in the literature analogous results to those in [35], [36], [15], [34] for the variable-coefficient Helmholtz equation.

1.2 Informal statement and discussion of the main results

The main results. This paper considers the variable-coefficient Helmholtz equation (1.4) with C^∞ coefficients posed in \mathbb{R}^d , $d = 2, 3$, with the Sommerfeld radiation condition at infinity. We obtain analogous results to those obtained in [35] for this scenario with constant coefficients. That is, we prove quasioptimality of the hp -FEM under the conditions (1.3) and provided that the solution operator is polynomially bounded in k ; see Theorem 3.4 below.

We obtain this result by decomposing the solution u to (1.4) into two components:

$$u|_{B_R} = u_{H^2} + u_{\mathcal{A}}$$

where $u_{H^2} \in H^2(B_R)$ and $u_{\mathcal{A}}$ is analytic in B_R , where B_R denotes the ball of radius R centred at the origin (and R is arbitrary); see Theorem 3.1 below. This is exactly analogous to the decomposition obtained in [35], except that now u satisfies the variable-coefficient equation (1.4) instead of (1.1).

Overview of the ideas behind the decomposition and subsequent bounds. The idea in [35] was to decompose the data f in (1.1) into “low-” and “high-” frequency components, with $u_{\mathcal{A}}$ the Helmholtz solution for the low-frequency component of f and u_{H^2} the Helmholtz solution for the high-frequency component of f . The frequency cut-offs were defining using the indicator function

$$1_{B_{\lambda k}}(\zeta) := \begin{cases} 1 & \text{for } |\zeta| \leq \lambda k, \\ 0 & \text{for } |\zeta| \geq \lambda k, \end{cases} \quad (1.5)$$

with λ a free parameter (see [35, Equation 3.31] and the surrounding text). In [35] the frequency cut-off (1.5) was then used with (a) the expression for u as a convolution of the fundamental solution and the data f , and (b) the fact that the fundamental solution is known explicitly when $A = I$ and $n = 1$, to obtain the appropriate bounds on $u_{\mathcal{A}}$ and u_{H^2} using explicit calculation.

In this paper we use the same idea as in [35] of decomposing into low- and high-frequency components, but apply frequency cut-offs to the solution u as opposed to the data f . Then, given any cut-off function that is zero for $|\zeta| \geq Ck$, bounding the corresponding low-frequency component $u_{\mathcal{A}}$ is relatively straightforward using basic properties of the Fourier-transform (namely the expression for the Fourier transform of a derivative and Parseval's theorem). Indeed, in Fourier space each derivative corresponds to a power of the Fourier variable ζ , and the frequency cut-off means that $|\zeta| \leq Ck$ for $u_{\mathcal{A}}$; i.e. every derivative of $u_{\mathcal{A}}$ brings down a power of k compared to $u_{\mathcal{A}}$ (see §5.3 below). The main difficulty therefore is showing that the high-frequency component u_{H^2} satisfies a bound with one power of k improvement over the bound satisfied by u .

The main idea of the present paper is that the high-frequency cut-off can be chosen so that the (scaled) Helmholtz operator

$$P_k := -(k^{-2}\nabla \cdot (A\nabla \cdot) + n) \quad (1.6)$$

is *semiclassically elliptic* on the support of the high-frequency cut-off. Furthermore, choosing the cut-off function to be smooth (as opposed to discontinuous, as in (1.5)) then allows us to use basic facts about the “nice” behaviour of elliptic semiclassical pseudodifferential operators (namely, they are invertible up to a small error) to prove the required bound on u_{H^2} . (Recall that semiclassical pseudodifferential operators are just pseudodifferential operators with a large/small parameter; in this case the large parameter is k .)

We now discuss further the frequency cut-offs and the bound on u_{H^2} via ellipticity.

The frequency cut-offs. In contrast to (1.5), we choose $\chi_\mu \in C_{\text{comp}}^\infty(\mathbb{R}^d)$ such that

$$\chi_\mu(k^{-2}|\zeta|^2) = \begin{cases} 1 & \text{for } |\zeta| \leq \sqrt{\mu}k, \\ 0 & \text{for } |\zeta| \geq \sqrt{2\mu}k, \end{cases} \quad (1.7)$$

where the parameter μ is chosen later in the argument. With the Fourier transform and its inverse defined by

$$\mathcal{F}\varphi(\zeta) := \int_{\mathbb{R}^d} \exp(-ix \cdot \zeta)\varphi(x) dx \quad \text{and} \quad \mathcal{F}^{-1}\psi(x) := (2\pi)^{-d} \int_{\mathbb{R}^d} \exp(ix \cdot \zeta)\psi(\zeta) d\zeta, \quad (1.8)$$

we define the low-frequency cut-off Π_L by

$$\Pi_L v(x) := \mathcal{F}^{-1}\left(\chi_\mu(k^{-2}|\zeta|^2)\mathcal{F}v(\zeta)\right), \quad (1.9)$$

and the high-frequency cut-off Π_H by

$$\Pi_H v(x) := \mathcal{F}^{-1}\left((1 - \chi_\mu(k^{-2}|\zeta|^2))\mathcal{F}v(\zeta)\right), \quad (1.10)$$

so that $\Pi_L + \Pi_H = I$. We let $\varphi \in C_c^\infty$ be equal to one on B_{R+1} and vanish outside B_{R+2} , and then

$$u_{\mathcal{A}} := \Pi_L(\varphi u)|_{B_R} \quad \text{and} \quad u_{H^2} := \Pi_H(\varphi u)|_{B_R}. \quad (1.11)$$

The bound on the high-frequency component u_{H^2} via ellipticity. Recall that a PDE is elliptic if its principal symbol is non-zero. The concept of ellipticity for semiclassical differential operators (or, more generally, semiclassical pseudodifferential operators) is analogous, except that it now involves the *semiclassical principal symbol* (see (4.17) below). The semiclassical principal symbol of P_k (1.6) is

$$\langle A\xi, \xi \rangle - n, \quad (1.12)$$

where $\langle \cdot, \cdot \rangle$ denotes the ℓ^2 inner product and $\xi = k^{-1}\zeta$ (see (4.12) below and the surrounding text).

If the parameter μ in the cut-off function χ_μ (1.7) is chosen to be a certain function of A and n (see (5.7) below), then the symbol (1.12) is bounded away from zero when $k^{-2}|\zeta|^2 \geq \mu$, i.e. in the

region of Fourier space where Π_H is non-zero; one therefore describes P_k as “microlocally elliptic”, where the adjective “microlocal” indicates that we have ellipticity on just a region of phase space (rather than on all of phase space in the more familiar global ellipticity).

These ellipticity properties are then used with the standard microlocal elliptic estimate for pseudodifferential operators, appearing in the semiclassical setting in, e.g., [14, Appendix E], and stated in this setting as Theorem 4.3 below. The whole point is that a semiclassical pseudodifferential operator that is elliptic in some region of phase space can be inverted (up to some small error) in that region, and the norm of the inverse is bounded uniformly in the large parameter (here k) as long as one uses weighted norms (analogous to the familiar H_k^1 norm (2.7)).

The result is that u_{H^2} satisfies a bound with one power of k improvement over the bound satisfied by u (compare (3.1) and (2.12)). To give a simple illustration of how ellipticity can give this improved k -dependence, we contrast the solutions of

$$P_k u := -(\Delta + k^2)u = f \quad \text{and} \quad \tilde{P}_k v := -(\Delta - k^2)v = f,$$

with both equations posed in \mathbb{R}^d with compactly-supported f , and with u satisfying the Sommerfeld radiation condition (1.2) and v satisfying boundedness at infinity. The $L^2 \rightarrow L^2$ bounds that are sharp in terms of k -dependence are

$$\|u\|_{L^2(B_R)} \lesssim k^{-1} \|f\|_{L^2(\mathbb{R}^d)} \quad \text{and} \quad \|v\|_{L^2(\mathbb{R}^d)} \lesssim k^{-2} \|f\|_{L^2(\mathbb{R}^d)},$$

with the former given by Part (i) of Theorem 2.7, and the latter following from the Lax-Milgram theorem. The operator P_k is not semiclassically elliptic on all of phase space (its semiclassical principal symbol is $|\xi|^2 - 1$), whereas \tilde{P}_k is semiclassically elliptic on all of phase space (its semiclassical principal symbol is $|\xi|^2 + 1$); we therefore see that ellipticity has resulted in the solution operator having improved k -dependence. The proof of the bound on u_{H^2} is more technical, but the idea – that the improvement in k -dependence comes from ellipticity – is the same.

The assumption that the solution operator is polynomially bounded in k . We need to assume that the solution operator is polynomially bounded in k (in sense of Definition 2.6 below), both in proving the bound on u_{H^2} , and in proving quasi-optimality of the hp -FEM.

The k -dependence of the Helmholtz solution operator depends on whether the problem is *trapping* or *nontrapping*. For the heterogeneous Helmholtz equation (1.4) posed in \mathbb{R}^d (i.e. with no obstacle), trapping can be created by the coefficients A and n ; see, e.g., [39]. If the problem is nontrapping, then the Helmholtz solution operator (measured in the natural norms) is bounded in k . However, under the strongest form of trapping, the Helmholtz solution operator can grow exponentially in k [39]. Nevertheless, it has recently been proved that, if a set of frequencies of arbitrarily small measure is excluded, then the solution operator is polynomially bounded under any type of trapping [28]. Therefore, the result that the hp -FEM is quasi-optimal holds for a wide class of Helmholtz problems; see Corollary 3.5 below.

Why do we need C^∞ coefficients? As highlighted above, our proof of the decomposition relies on standard results about semiclassical pseudodifferential operators (recapped in §4). These results are usually stated for C^∞ symbols, and thus to fit into this framework A and n must be C^∞ . However, examining the results we use, we see that we only need the symbol of the PDE to be in C^L where L depends only on the dimension d and on the exponent M appearing in the assumption that the solution operator is polynomially bounded (see Definitions 2.5 and 2.6 below). Therefore, while we consider $A, n \in C^\infty$ to easily use results about semiclassical pseudodifferential operators from [52], [14, Appendix E], our results hold for $A \in C^L$ and $n \in C^L$, where $L = L(d, M)$.

Extending the decomposition result to the solution of other PDEs. Our proof of the decomposition result only relies on the principal symbol of the differential operator being bounded below at infinity (in the sense of (3.8) below). Therefore, the decomposition result Theorem 3.1 is valid for a much larger class of PDEs (and indeed pseudodifferential operators) than (1.4); see Remark 3.7 below for more details.

In the follow-up paper [27], we use the ideas of the present paper combined with much more sophisticated tools of semiclassical and microlocal analysis (namely the black-box scattering framework of Sjöstrand–Zworski [45], the Helffer–Sjöstrand functional calculus [23], and associated results by Helffer, Robert, and Sjöstrand [22], [40], [44]) to prove analogous decompositions for a wide variety of scattering problems (albeit with slightly weaker estimates on $u_{\mathcal{A}}$). In particular, the main result of the present paper, Theorem 3.1, is rederived in this more general context as [27, Theorem 1.16].

We also note that, as announced in the abstract [4], Bernkopf, Chaumont–Frelet, and Melenk are also studying the question of k -explicit convergence of the hp -FEM for the Helmholtz equation with variable coefficients.

Outline of the paper. §2 gives the definitions of the boundary-value problem and the finite-element method. §3 states the main results. §4 recaps results about semiclassical pseudodifferential operators, with [52] and [14, Appendix E] as the main references. §5 proves the result about the decomposition $u|_{B_R} = u_{H^2} + u_{\mathcal{A}}$ (Theorem 3.1). §6 proves the result about quasioptimality of the hp -FEM (Theorem 3.4).

2 Formulation of the problem

2.1 The boundary value problem

Assumption 2.1 (Assumptions on the coefficients) $\mathbf{A} \in C^\infty(\mathbb{R}^d, \text{SPD})$ (where SPD is the set of $d \times d$ real, symmetric, positive-definite matrices) is such that $\text{supp}(1 - \mathbf{A})$ is compact in \mathbb{R}^d and there exist $0 < A_{\min} \leq A_{\max} < \infty$ such that, in the sense of quadratic forms,

$$A_{\min} \leq \mathbf{A}(x) \leq A_{\max} \quad \text{for all } x \in \mathbb{R}^d. \quad (2.1)$$

$n \in C^\infty(\mathbb{R}^d, \mathbb{R})$ is such that $\text{supp}(1 - n)$ is compact in \mathbb{R}^d and there exist $0 < n_{\min} \leq n_{\max} < \infty$ such that

$$n_{\min} \leq n(x) \leq n_{\max} \quad \text{for all } x \in \mathbb{R}^d. \quad (2.2)$$

Let $R > 0$ be such that $\text{supp}(1 - \mathbf{A}) \cup \text{supp}(1 - n) \Subset B_R$, where B_R denotes the ball of radius R about the origin and \Subset denotes compact containment. Let γ and ∂_n denote the Dirichlet and Neumann traces, respectively, on ∂B_R , where the normal vector for the Neumann trace points out of B_R .

Define $\text{DtN}_k : H^{1/2}(\partial B_R) \rightarrow H^{-1/2}(\partial B_R)$ to be the Dirichlet-to-Neumann map for the equation $\Delta u + k^2 u = 0$ posed in the exterior of B_R with the Sommerfeld radiation condition (1.2). The definition of DtN_k in terms of Hankel functions and polar coordinates (when $d = 2$)/spherical polar coordinates (when $d = 3$) is given in, e.g., [35, Equations 3.7 and 3.10].

Definition 2.2 (Heterogeneous Helmholtz Problem on \mathbb{R}^d) Given \mathbf{A} and n satisfying Assumption 2.1, $R > 0$ such that $\text{supp}(1 - \mathbf{A}) \cup \text{supp}(1 - n) \Subset B_R$, $k > 0$, and $F \in (H^1(B_R))^*$, $u \in H^1(B_R)$ satisfies the Heterogeneous Helmholtz Problem on \mathbb{R}^d if u satisfies the variational problem

$$\text{find } u \in H^1(B_R) \text{ such that } a(u, v) = F(v) \quad \text{for all } v \in H^1(B_R), \quad (2.3)$$

where

$$a(u, v) := \int_{B_R} \left((\mathbf{A} \nabla u) \cdot \overline{\nabla v} - k^2 n u \bar{v} \right) - \langle \text{DtN}_k(\gamma u), \gamma v \rangle_{\partial B_R}, \quad (2.4)$$

where $\langle \cdot, \cdot \rangle_{\partial B_R}$ denotes the duality pairing on ∂B_R that is linear in the first argument and antilinear in the second.

Lemma 2.3 (Helmholtz boundary value problems included in Definition 2.2)

(i) If

$$F(v) := \int_{B_R} f \bar{v} \quad (2.5)$$

with $f \in L^2(B_R)$, then the solution u to (2.3) equals $\tilde{u}|_{B_R}$, where $\tilde{u} \in H_{\text{loc}}^1(\mathbb{R}^d)$ is the solution to

$$\nabla \cdot (\mathbf{A}\nabla\tilde{u}) + k^2 n\tilde{u} = -f \quad \text{in } \mathbb{R}^d,$$

and \tilde{u} satisfies the Sommerfeld radiation condition (1.2).

(ii) If

$$F(v) := \int_{\partial B_R} (\partial_n u^I - \text{DtN}_k(\gamma u^I)) \bar{\gamma} v \quad \text{with} \quad u^I(x) := \exp(ikx \cdot a), \quad (2.6)$$

where $a \in \mathbb{R}^d$ with $|a| = 1$, then the solution u to (2.3) equals $\tilde{u}|_{B_R}$, where $\tilde{u} \in H_{\text{loc}}^1(\mathbb{R}^d)$ is the solution of the Helmholtz plane-wave scattering problem; i.e.

$$\nabla \cdot (\mathbf{A}\nabla\tilde{u}) + k^2 n\tilde{u} = 0 \quad \text{in } \mathbb{R}^d,$$

and $\tilde{u}^S := \tilde{u} - u^I$ satisfies the Sommerfeld radiation condition (1.2).

Part (i) of Lemma 2.3 is proved in, e.g., [20, Lemma 3.3]; the proof of Part (ii) is similar.

Let the weighted H^1 norm, $\|\cdot\|_{H_k^1(B_R)}$, be defined by

$$\|u\|_{H_k^1(B_R)}^2 := \|\nabla u\|_{L^2(B_R)}^2 + k^2 \|u\|_{L^2(B_R)}^2. \quad (2.7)$$

Lemma 2.4 *The solution of the Heterogeneous Helmholtz Problem on \mathbb{R}^d (defined in Definition 2.2) exists, is unique, and there exists $C(k, \mathbf{A}, n, R) > 0$ such that*

$$\|u\|_{H_k^1(B_R)} \leq C \|F\|_{(H_k^1(B_R))^*} \quad \text{for all } k > 0. \quad (2.8)$$

Proof. Uniqueness follows from the unique continuation principle; see [20, §1], [21, §2] and the references therein. Since $a(\cdot, \cdot)$ satisfies a Gårding inequality (see (6.4) below), Fredholm theory then gives existence and the bound (2.8). \blacksquare

Properties of DtN_k and $a(\cdot, \cdot)$. We use later the following two properties of DtN_k : given $k_0, R_0 > 0$, there exists $C_{\text{DtN}} = C_{\text{DtN}}(k_0, R_0)$ such that, for all $k \geq k_0$ and $R \geq R_0$,

$$|\langle \text{DtN}_k(\gamma u), \gamma v \rangle_{\partial B_R}| \leq C_{\text{DtN}1} \|u\|_{H_k^1(B_R)} \|v\|_{H_k^1(B_R)} \quad (2.9)$$

for all $u, v \in H^1(B_R)$, and

$$-\Re \langle \text{DtN}_k \phi, \phi \rangle_{\partial B_R} \geq 0 \quad \text{for all } \phi \in H^{1/2}(\partial B_R). \quad (2.10)$$

For a proof of (2.9), see [35, Lemma 3.3]. For a proof of (2.10), see [37, Theorem 2.6.4] (for $d = 3$) and [7, Corollary 3.1] or [35, Lemma 3.10] (for $d = 2, 3$).

Let $C_{\text{cont}} = C_{\text{cont}}(\mathbf{A}, n, R, k_0)$ be the *continuity constant* of the sesquilinear form $a(\cdot, \cdot)$ (defined in (2.4)) in the norm $\|\cdot\|_{H_k^1(B_R)}$; i.e.

$$|a(u, v)| \leq C_{\text{cont}} \|u\|_{H_k^1(B_R)} \|v\|_{H_k^1(B_R)} \quad \text{for all } u, v \in H^1(B_R) \text{ and } k \geq k_0.$$

By the Cauchy-Schwarz inequality and (2.9),

$$C_{\text{cont}} \leq \max\{A_{\text{max}}, n_{\text{max}}\} + C_{\text{DtN}1}. \quad (2.11)$$

2.2 The behaviour of the solution operator for large k

Definition 2.5 (C_{sol}) *Given $f \in L^2(B_R)$, let u be the solution of the heterogeneous Helmholtz equation (1.4) with the Sommerfeld radiation condition (1.2) (i.e. u is the solution of the variational problem (2.3) with $F(v)$ given by (2.5)). Given $k_0 > 0$, let $C_{\text{sol}} = C_{\text{sol}}(k, \mathbf{A}, n, R, k_0) > 0$ be such that*

$$\|u\|_{H_k^1(B_R)} \leq C_{\text{sol}} \|f\|_{L^2(B_R)} \quad \text{for all } k > 0. \quad (2.12)$$

C_{sol} exists by Lemma 2.4; indeed, with C given by (2.8), $C_{\text{sol}} := C/k$.

How C_{sol} depends on k is crucial to the analysis below, and to emphasise this we write $C_{\text{sol}} = C_{\text{sol}}(k)$. Below we consider C_{sol} with different values of R , and we then write, e.g., $C_{\text{sol}}(k; R)$ (as in the bound (3.2) below).

A key assumption in the analysis of the Helmholtz hp -FEM is that $C_{\text{sol}}(k)$ is polynomially bounded in k in the following sense.

Definition 2.6 (C_{sol} is polynomially bounded in k) *Given k_0 and $K \subset [k_0, \infty)$, $C_{\text{sol}}(k)$ is polynomially bounded for $k \in K$ if there exists $C > 0$ and $M > 0$ such that*

$$C_{\text{sol}}(k) \leq Ck^M \text{ for all } k \in K, \quad (2.13)$$

where C and M are independent of k (but depend on k_0 and possibly also on K, A, n, d, R).

There exist C^∞ coefficients A and n such that $C_{\text{sol}}(k_j) \geq c_1 \exp(c_2 k_j)$ for $0 < k_1 < k_2 < \dots$ with $k_j \rightarrow \infty$ as $j \rightarrow \infty$, see [39], but this exponential growth is the worst-possible, since $C_{\text{sol}}(k) \leq c_3 \exp(c_4 k)$ for all $k \geq k_0$ by [5, Theorem 2]. We now recall results on when $C_{\text{sol}}(k)$ is polynomially bounded in k .

Theorem 2.7 (Conditions under which $C_{\text{sol}}(k)$ is polynomially bounded in k)

(i) A and n are C^∞ and nontrapping (i.e. all the trajectories of the Hamiltonian flow defined by the symbol of (1.4) starting in B_R leave B_R after a uniform time), then $C_{\text{sol}}(k)$ is independent of k for all k , i.e., (2.13) holds for all k with $M = 0$.

(ii) If $n = 1$ and A is $C^{0,1}$ then, given $k_0 > 0$ and $\delta > 0$ there exists a set $J \subset [k_0, \infty)$ with $|J| \leq \delta$ such that

$$C_{\text{sol}}(k) \leq Ck^{5d/2+1+\varepsilon} \text{ for all } k \in [k_0, \infty) \setminus J, \quad (2.14)$$

for any $\varepsilon > 0$, where C depends on $\delta, \varepsilon, d, k_0$, and A . If A is $C^{1,\sigma}$ for some $\sigma > 0$ then the exponent is reduced to $5d/2 + \varepsilon$.

References for the proof.

(i) is proved using either (a) the propagation of singularities results of [13] combined with either the parametrix argument of [48, Theorem 3]/ [49, Chapter 10, Theorem 2] or Lax–Phillips theory [30], or (b) the defect-measure argument of [6, Theorem 1.3 and §3]. It has recently been proved that, for this situation, C_{sol} is proportional to the length of the longest trajectory in B_R ; see [16, Theorems 1 and 2, and Equation 6.32].

(ii) is proved in [28, Theorem 1.1 and Corollary 3.6]. ■

2.3 The finite-element method

Let $(V_N)_{N=0}^\infty$ be a sequence of finite-dimensional subspaces of $H^1(B_R)$ that converge to $H^1(B_R)$ in the sense that, for all $v \in H^1(B_R)$,

$$\lim_{N \rightarrow \infty} \left(\min_{v_N \in V_N} \|v - v_N\|_{H^1(B_R)} \right) = 0.$$

Later we specialise to the triangulations described in [35, §5], which allow curved elements and thus fit ∂B_R exactly.

The finite-element method for the variational problem (2.3) is the Galerkin method applied to the variational problem (2.3), i.e.

$$\text{find } u_N \in V_N \text{ such that } a(u_N, v_N) = F(v_N) \text{ for all } v_N \in V_N. \quad (2.15)$$

3 Statement of the main results

Theorem 3.1 (Decomposition of the solution) *Let A and n satisfy Assumption 2.1 and let $R > 0$ be such that $\text{supp}(1 - A) \cup \text{supp}(1 - n) \Subset B_R$. Given $f \in L^2(B_R)$, let u satisfy $\nabla \cdot (A \nabla u) + k^2 n u = -f$ in \mathbb{R}^d and the Sommerfeld radiation condition (1.2).*

If $C_{\text{sol}}(k)$ is polynomially bounded (in the sense of Definition 2.6) for $k \in K \subset [k_0, \infty)$, then there exist $C_3, C_4, C_5 > 0$ such that

$$u|_{B_R} = u_{H^2} + u_{\mathcal{A}}$$

where $u_{H^2} \in H^2(B_R)$ with

$$\|\partial^\alpha u_{H^2}\|_{L^2(B_R)} \leq C_3 k^{|\alpha|-2} \|f\|_{L^2(B_R)} \quad \text{for all } |\alpha| \leq 2 \text{ and for all } k \in K \subset [k_0, \infty), \quad (3.1)$$

and $u_{\mathcal{A}} \in C^\infty(B_R)$ with

$$\|\partial^\beta u_{\mathcal{A}}\|_{L^2(B_R)} \leq C_{\text{sol}}(k; R+2) C_4 (C_5 k)^{|\beta|-1} \|f\|_{L^2(B_R)} \quad \text{for all } \beta \text{ and for all } k \in K \subset [k_0, \infty), \quad (3.2)$$

where C_3, C_4 , and C_5 depend on \mathbf{A}, n, d , and k_0 , but are independent of k, f, α , and β .

Remark 3.2 ($u_{\mathcal{A}}$ is analytic) Since C_4 and C_5 are independent of β , the bound (3.2) implies that $u_{\mathcal{A}}$ is in the class of analytic functions on B_R , $\mathcal{A}(B_R)$, defined by

$$\mathcal{A}(B_R) := \left\{ v \in \bigcap_{n \in \mathbb{N}} H^n(B_R) : \exists c_0, c_1 > 0, \text{ independent of } n, \text{ such that } |u|_{H^n(B_R)} \leq c_1 c_0^n n! \right\},$$

where $|u|_{H^n}^2 := \sum_{|\alpha|=n} \|\partial^\alpha u\|_{L^2}^2$. See, e.g., [11, §1.1.b], both for this definition, and for how the definition implies convergence of the Taylor series of elements of $\mathcal{A}(B_R)$ at every point in $\overline{B_R}$.

Remark 3.3 (The bounds of Theorem 3.1 written with the notation ∇^n) The analogous bounds to (3.1) and (3.2) in [35], [36] are written using the notation

$$|\nabla^n u(x)|^2 := \sum_{|\alpha|=n} \frac{n!}{\alpha!} |\partial^\alpha u(x)|^2.$$

Since $\sum_{|\alpha|=n} (n!/\alpha!) = d^n$,

$$\text{if } \|\partial^\alpha u\|_{L^2(B_R)} \leq C_1 (C_2)^{|\alpha|} \text{ for all } \alpha \text{ with } |\alpha| = n, \text{ then } \|\nabla^n u\|_{L^2(B_R)} \leq C_1 (C_2 \sqrt{d})^n,$$

and so the bounds (3.1) and (3.2) can also be written as bounds on $\|\nabla^n u_{H^2}\|_{L^2(B_R)}$ and $\|\nabla^n u_{\mathcal{A}}\|_{L^2(B_R)}$ respectively.

The following result about quasioptimality of the hp -FEM is then obtained by combining Theorem 3.1, well-known results about the convergence of the Galerkin method based on duality arguments (recapped in Lemma 6.4 below), and results about the hp approximation spaces in [35, §5] (used in Lemma 6.5 below).

Theorem 3.4 (Quasioptimality of the hp -FEM if $C_{\text{sol}}(k)$ is polynomially bounded) Let $d = 2$ or 3 , and let $k_0 > 0$. Let $(V_N)_{N=0}^\infty$ be the piecewise-polynomial approximation spaces described in [35, §5] (where, in particular, the triangulations are quasi-uniform), and let u_N be the Galerkin solution defined by (2.15).

If $C_{\text{sol}}(k)$ is polynomially bounded (in the sense of Definition 2.6) for $k \in K \subset [k_0, \infty)$ then there exist $C_1, C_2 > 0$, depending on \mathbf{A}, n, R , and d , and k_0 , but independent of k, h , and p , such that if (1.3) holds, then, for all $k \in K$, the Galerkin solution exists, is unique, and satisfies the quasi-optimal error bound

$$\|u - u_N\|_{H_k^1(B_R)} \leq C_{\text{qo}} \min_{v_N \in V_N} \|u - v_N\|_{H_k^1(B_R)}, \quad (3.3)$$

with

$$C_{\text{qo}} := \frac{2(\max\{A_{\max}, n_{\max}\} + C_{\text{DtN1}})}{A_{\min}} \quad (3.4)$$

Combining Theorem 3.4 with the results on $C_{\text{sol}}(k)$ recapped in Theorem 2.7, we obtain the following specific examples of coefficients \mathbf{A} and n when quasioptimality holds.

Corollary 3.5 (Quasioptimality under specific conditions on A and n) *Let $d = 2$ or 3 , and let $k_0 > 0$.*

(i) If A and n are nontrapping, then there exist $C_1, C_2 > 0$, depending on A, n, R , and d , and k_0 , but independent of k, h , and p , such that if (1.3) holds then, for all $k \geq k_0$, the Galerkin solution exists, is unique, and satisfies the quasi-optimal error bound (3.3) with C_{qo} given by (3.4).

(ii) If A is C^∞ and $n = 1$ then, given $\delta > 0$, there exist a set J with $|J| \leq \delta$ and constants $\tilde{C}_1, \tilde{C}_2 > 0$, with all three depending on A, n, R, d , and k_0 , but independent of k , and \tilde{C}_2 additionally depending on δ and k_0 such that, for all $k \in [k_0, \infty) \setminus J$, if (1.3) holds (with C_1, C_2 replaced by \tilde{C}_1, \tilde{C}_2) then the Galerkin solution exists, is unique, and satisfies (3.3) with C_{qo} given by (3.4).

For the plane-wave scattering problem (i.e. for $F(v)$ given by (2.6)), the regularity result

$$|u|_{H^2(B_R)} \leq C_{\text{osc}} k \|u\|_{H_k^1(B_R)} \quad (3.5)$$

was recently proved in [29, Theorem 9.1 and Remark 9.10], where C_{osc} depends on A, n, d , and R , but is independent of k . The polynomial approximation bounds in [35, §B] imply that, for the sequence of approximation spaces $(V_N)_{N=0}^\infty$ described in [35, §5],

$$\min_{v_N \in V_N} \|u - v_N\|_{H_k^1(B_R)} \leq C_6 \frac{h}{p} \left(1 + \frac{kh}{p}\right) |u|_{H^2(B_R)} \quad (3.6)$$

where C_6 only depends on the constants in [35, Assumption 5.2] (which depend on the element maps from the reference element). Using (3.6) and (3.5) to bound the right-hand side of (3.3), we obtain the following bound on the relative error of the Galerkin solution.

Corollary 3.6 (Bound on the relative error of the Galerkin solution) *Let the assumptions of Theorem 3.4 hold and, furthermore, let $F(v)$ be given by (2.6) (so that u is the solution of the plane-wave scattering problem). If $C_{\text{sol}}(k)$ is polynomially bounded (in the sense of Definition 2.6) for $k \in K \subset [k_0, \infty)$, then there exists $C_6 > 0$, independent of k, h , and p , such that if (1.3) holds, then, for all $k \in K$,*

$$\frac{\|u - u_N\|_{H_k^1(B_R)}}{\|u\|_{H_k^1(B_R)}} \leq C_{qo} C_6 C_{\text{osc}} C_1 (1 + C_1), \quad (3.7)$$

with C_{qo} given by (3.4); i.e. the relative error can be made arbitrarily small by making C_1 smaller.

Remark 3.7 (Theorem 3.1 is valid for solutions of a much larger class of PDEs)

Inspecting the proof of Theorem 3.1 below, we see that the conclusion, i.e. the decomposition $u = u_{H^2} + u_A$ with u_{H^2} and u_A satisfying the bounds (3.1) and (3.2) respectively, holds under much weaker assumptions. Indeed, the conclusion still holds under the following three assumptions only.

(i) P_k is a family of properly-supported second-order pseudo-differential operators, with principal symbol $p_k(x, \zeta)$,

(ii) $p_k(x, \zeta)$ is coercive at infinity in the sense that

$$\liminf_{|\xi| \rightarrow \infty, x \in \mathbb{R}^d} \langle k\xi \rangle^{-2} p_k(x, k\xi) \geq c > 0, \quad (3.8)$$

where $c > 0$ does not depend on k , and

(iii) the solution to $P_k u = -f$, posed in \mathbb{R}^d with $\text{supp } f \subset B_R$ and $f \in L^2(B_R)$, satisfies the bound

$$\|u\|_{L^2(B_{R+2})} \leq C k^M \|f\|_{L^2(B_R)},$$

with C and M independent of k, u , and f . (In fact, the 2 in the $R+2$ on the left-hand side of the bound can be replaced by any number > 0 .)

In particular, no assumption is made about lower-order terms of P_k , or the behaviour of u at infinity (such as a radiation condition).

4 Recap of relevant results about semiclassical pseudodifferential operators

The proof of Theorem 3.1 relies on standard results about semiclassical pseudodifferential operators. We review these here, with our default references being [52] and [14, Appendix E]. Homogeneous – as opposed to semiclassical – versions of the results in this section can be found in, e.g., [47, Chapter 7], [41, Chapter 7], [25, Chapter 6].¹

While the use of homogeneous pseudodifferential operators in numerical analysis is well established, see, e.g., [41], [25], there has been less use of semiclassical pseudodifferential operators. However, these are ideally-suited for studying the high-frequency behaviour of Helmholtz solutions. Indeed, semiclassical pseudodifferential operators are just pseudodifferential operators with a large/small parameter, and behaviour with respect to this parameter is then explicitly kept track of in the associated calculus.

The semiclassical parameter $\hbar = k^{-1}$. Instead of working with the parameter k and being interested in the large- k limit, the semiclassical literature usually works with a parameter $\hbar := k^{-1}$ and is interested in the small- \hbar limit. So that we can easily recall results from this literature, we also work with the small parameter k^{-1} , but to avoid a notational clash with the meshwidth of the FEM, we let $\hbar := k^{-1}$ (the notation \hbar comes from the fact that the semiclassical parameter is related to Planck’s constant, which is written as $2\pi\hbar$; see, e.g., [52, §1.2], [14, Page 82], [32, Chapter 1]). In this notation, the Helmholtz equation $\nabla \cdot (A\nabla u) + k^2 nu = -f$ becomes

$$P_\hbar u = \hbar^2 f, \quad \text{where} \quad P_\hbar := -\hbar^2 \nabla \cdot (A\nabla \cdot) - n. \quad (4.1)$$

While some results in semiclassical analysis are valid in the limit \hbar small, the results we recap in this section are valid for all $0 < \hbar \leq \hbar_0$ with $\hbar_0 < \infty$ arbitrary.

The semiclassical Fourier transform \mathcal{F}_\hbar . The semiclassical Fourier transform is defined for $\hbar > 0$ by

$$\mathcal{F}_\hbar \phi(\xi) := \int_{\mathbb{R}^d} \exp(-ix \cdot \xi/\hbar) \phi(x) \, dx,$$

and its inverse by

$$\mathcal{F}_\hbar^{-1} \psi(x) := (2\pi\hbar)^{-d} \int_{\mathbb{R}^d} \exp(ix \cdot \xi/\hbar) \psi(\xi) \, d\xi; \quad (4.2)$$

see [52, §3.3]. Then

$$\mathcal{F}_\hbar \left((-i\hbar\partial)^\alpha \phi \right) = \xi^\alpha \mathcal{F}_\hbar \phi \quad (4.3)$$

and

$$\|\phi\|_{L^2(\mathbb{R}^d)} = \frac{1}{(2\pi\hbar)^{d/2}} \|\mathcal{F}_\hbar \phi\|_{L^2(\mathbb{R}^d)}. \quad (4.4)$$

Semiclassical Sobolev spaces. In the same way that it is convenient to work with the weighted H^1 norm (2.7) when studying the Helmholtz equation with parameter k , it is convenient to use norms weighted with \hbar when studying (4.1). Therefore on the space

$$H_\hbar^s(\mathbb{R}^d) := \left\{ u \in L^2(\mathbb{R}^d), \langle \xi \rangle^s \mathcal{F}_\hbar u \in L^2(\mathbb{R}^d) \right\}, \quad \text{where } \langle \xi \rangle := (1 + |\xi|^2)^{1/2}, \quad s \in \mathbb{R},$$

we use the norm

$$\|u\|_{H_\hbar^s(\mathbb{R}^d)}^2 := (2\pi\hbar)^{-d} \int_{\mathbb{R}^d} \langle \xi \rangle^{2s} |\mathcal{F}_\hbar u(\xi)|^2 \, d\xi; \quad (4.5)$$

see [52, §8.3], [14, §E.1.8]. We abbreviate $H_\hbar^s(\mathbb{R}^d)$ to H_\hbar^s and $L^2(\mathbb{R}^d)$ to L^2 .

¹The counterpart of “semiclassical” involving differential/pseudodifferential operators without a small parameter is usually called “homogeneous” (owing to the homogeneity of the principal symbol) rather than “classical.” “Classical” describes the behaviour in either calculus in the small- \hbar or high-frequency limit respectively, where commutators of operators become Poisson brackets of symbols, hence classical particle dynamics replaces wave motion.

We record for later the fact that, by (4.3) and (4.4), for multiindices α ,

$$\hbar^{|\alpha|} \|\partial^\alpha \phi\|_{L^2} = \|(-i\hbar\partial)^\alpha \phi\|_{L^2} = \frac{1}{(2\pi\hbar)^{d/2}} \|\xi^\alpha \mathcal{F}_\hbar \phi\|_{L^2} \leq \frac{1}{(2\pi\hbar)^{d/2}} \|\langle \xi \rangle^{|\alpha|} \mathcal{F}_\hbar \phi\|_{L^2} = \|\phi\|_{H_\hbar^{|\alpha|}}. \quad (4.6)$$

Phase space. The set of all possible positions x and momenta (i.e. Fourier variables) ξ is denoted by $T^*\mathbb{R}^d$; this is known informally as “phase space”. Strictly, $T^*\mathbb{R}^d := \mathbb{R}^d \times (\mathbb{R}^d)^*$, but for our purposes, we can consider $T^*\mathbb{R}^d$ as $\{(x, \xi) : x \in \mathbb{R}^d, \xi \in \mathbb{R}^d\}$.

To deal with the behavior of functions on phase space uniformly near $\xi = \infty$ (so-called *fiber infinity*), we consider the *radial compactification* in the ξ variable of $T^*\mathbb{R}^d$. This is defined by

$$\overline{T^*\mathbb{R}^d} := \mathbb{R}^d \times B^d,$$

where B^d denotes the closed unit ball, considered as the closure of the image of \mathbb{R}^d under the radial compactification map

$$\text{RC} : \xi \mapsto \xi/(1 + \langle \xi \rangle);$$

see [14, §E.1.3]. Near the boundary of the ball, $|\xi|^{-1} \circ \text{RC}^{-1}$ is a smooth function, vanishing to first order at the boundary, with $(|\xi|^{-1} \circ \text{RC}^{-1}, \widehat{\xi} \circ \text{RC}^{-1})$ thus giving local coordinates on the ball near its boundary. The boundary of the ball should be considered as a sphere at infinity consisting of all possible *directions* of the momentum variable. More generally, we denote $\overline{T^*X} := X \times B^d$ for $X \subset \mathbb{R}^d$, and where appropriate (e.g., in dealing with finite values of ξ only), we abuse notation by dropping the composition with RC from our notation and simply identifying \mathbb{R}^d with the interior of B^d .

Symbols, quantisation, and semiclassical pseudodifferential operators. A symbol is a function on $T^*\mathbb{R}^d$ that is also allowed to depend on \hbar , and thus can be considered as an \hbar -dependent family of functions. Such a family $a = (a_\hbar)_{0 < \hbar \leq \hbar_0}$, with $a_\hbar \in C^\infty(T^*\mathbb{R}^d)$, is a *symbol of order m* , written as $a \in S^m(\mathbb{R}^d)$, if for any multiindices α, β

$$|\partial_x^\alpha \partial_\xi^\beta a(x, \xi)| \leq C_{\alpha, \beta} \langle \xi \rangle^{m - |\beta|} \quad \text{for all } (x, \xi) \in T^*\mathbb{R}^d \text{ and for all } 0 < \hbar \leq \hbar_0, \quad (4.7)$$

where $C_{\alpha, \beta}$ does not depend on \hbar , x , or ξ ; see [52, p. 207], [14, §E.1.2]. In this paper, we only consider these symbol classes on \mathbb{R}^d , and so we abbreviate $S^m(\mathbb{R}^d)$ to S^m .

For $a \in S^m$, we define the *semiclassical quantisation* of a , $\text{Op}_\hbar(a) : \mathcal{S}(\mathbb{R}^d) \rightarrow \mathcal{S}(\mathbb{R}^d)$, by

$$(\text{Op}_\hbar(a)v)(x) := (2\pi\hbar)^{-d} \int_{\mathbb{R}^d} \int_{\mathbb{R}^d} \exp(i(x-y) \cdot \xi/\hbar) a(x, \xi) v(y) dy d\xi \quad (4.8)$$

for $v \in \mathcal{S}(\mathbb{R}^d)$; [52, §4.1] [14, Page 543]. The integral in (4.8) need not converge, and can be understood *either* as an oscillatory integral in the sense of [52, §3.6], [24, §7.8], *or* as an iterated integral, with the y integration performed first; see [14, Page 543].

Conversely, if A can be written in the form above, i.e. $A = \text{Op}_\hbar(a)$ with $a \in S^m$, we say that A is a *semiclassical pseudo-differential operator of order m* and we write $A \in \Psi_\hbar^m$. We use the notation $a \in \hbar^l S^m$ if $\hbar^{-l} a \in S^m$; similarly $A \in \hbar^l \Psi_\hbar^m$ if $\hbar^{-l} A \in \Psi_\hbar^m$.

Theorem 4.1 (Composition and mapping properties of semiclassical pseudo-differential operators [52, Theorem 8.10], [14, Proposition E.17 and Proposition E.19].)

If $A \in \Psi_\hbar^{m_1}$ and $B \in \Psi_\hbar^{m_2}$, then

$$(i) \quad AB \in \Psi_\hbar^{m_1+m_2},$$

$$(ii) \quad [A, B] := AB - BA \in \hbar \Psi_\hbar^{m_1+m_2-1},$$

(iii) For any $s \in \mathbb{R}$, A is bounded uniformly in \hbar as an operator from H_\hbar^s to $H_\hbar^{s-m_1}$.

Residual class. We say that $A = O(\hbar^\infty)_{\Psi^{-\infty}}$ if, for any $s > 0$ and $N \geq 1$, there exists $C_{s,N} > 0$ so that

$$\|A\|_{H_h^{-s} \rightarrow H_h^s} \leq C_{N,s} \hbar^N; \quad (4.9)$$

i.e. $A \in \Psi_h^{-\infty}$ and furthermore all of its operator norms are bounded by any algebraic power of \hbar .

Principal symbol σ_h . Let the quotient space $S^m / \hbar S^{m-1}$ be defined by identifying elements of S^m that differ only by an element of $\hbar S^{m-1}$. For any m , there is a linear, surjective map

$$\sigma_h^m : \Psi_h^m \rightarrow S^m / \hbar S^{m-1},$$

called the *principal symbol map*, such that, for $a \in S^m$,

$$\sigma_h^m(\text{Op}_h(a)) = a \pmod{\hbar S^{m-1}}; \quad (4.10)$$

see [52, Page 213], [14, Proposition E.14] (observe that (4.10) implies that $\ker(\sigma_h^m) = \hbar \Psi_h^{m-1}$).

When applying the map σ_h^m to elements of Ψ_h^m , we denote it by σ_h (i.e. we omit the m dependence) and we use $\sigma_h(A)$ to denote one of the representatives in S^m (with the results we use then independent of the choice of representative). Key properties of the principal symbol that we use below are that

$$\sigma_h(AB) = \sigma_h(A)\sigma_h(B), \quad (4.11)$$

$$\sigma_h(P_h) = \langle A\xi, \xi \rangle - n, \quad (4.12)$$

where $\langle \cdot, \cdot \rangle$ denotes the ℓ^2 inner product on \mathbb{R}^d . The property (4.11) is proved in [14, Proposition E.17], (4.12) follows from (4.10) since $P_h = \text{Op}_h(\langle A\xi, \xi \rangle - n - i\hbar\xi_\ell \partial_j A_{j\ell})$ (where we sum over the indices j and ℓ).

Operator wavefront set WF_h . We say that $(x_0, \xi_0) \in \overline{T^* \mathbb{R}^d}$ is *not* in the *semiclassical operator wavefront set* of $A = \text{Op}_h(a) \in \Psi_h^m$, denoted by $\text{WF}_h A$, if there exists a neighbourhood U of (x_0, ξ_0) such that for all multiindices α, β and all $N \geq 1$ there exists $C_{\alpha,\beta,U,N} > 0$ (independent of \hbar) so that, for all $0 < \hbar \leq \hbar_0$,

$$|\partial_x^\alpha \partial_\xi^\beta a(x, \xi)| \leq C_{\alpha,\beta,U,N} \hbar^N \langle \xi \rangle^{-N} \quad \text{for all } (x, \text{RC}(\xi)) \in U; \quad (4.13)$$

i.e. outside its semiclassical operator wavefront set an operator vanishes faster than any algebraic power of both \hbar and $\langle \xi \rangle^{-1}$; see [52, Page 194], [14, Definition E.27]. Three properties of the semiclassical operator wavefront set that we use below are

$$\text{WF}_h(AB) \subset \text{WF}_h A \cap \text{WF}_h B \quad (4.14)$$

(see [52, §8.4], [14, E.2.5]),

$$\text{WF}_h(\text{Op}_h(a)) \subset \text{supp } a \quad (4.15)$$

(since $(\text{supp } a)^c \subset (\text{WF}_h(\text{Op}_h(a)))^c$ by (4.13)), and

$$\text{WF}_h A = \emptyset \iff A = O(\hbar^\infty)_{\Psi^{-\infty}} \quad (4.16)$$

(see [14, E.2.2]).

Compactly-supported operators. We say that A is *compactly supported* if its Schwartz kernel is compactly supported in some set $K \Subset \mathbb{R}^d \times \mathbb{R}^d$, for all $0 < \hbar \leq \hbar_0$. We recall that if $\mathcal{D}(\mathbb{R}^d) := C_{\text{comp}}^\infty(\mathbb{R}^d)$ (i.e. the set of test functions) and $\mathcal{D}'(\mathbb{R}^d)$ denote the set of linear functionals on $\mathcal{D}(\mathbb{R}^d)$ (i.e. the set of distributions), given a bounded, sequentially-continuous operator $A : \mathcal{D} \rightarrow \mathcal{D}'$ there exists a *Schwartz kernel* $\mathcal{K}_A \in \mathcal{D}'(\mathbb{R}^d \times \mathbb{R}^d)$ such that

$$Av(x) = \int_{\mathbb{R}^d} \mathcal{K}_A(x, y)v(y) dy,$$

in the sense of distributions; see, e.g., [24, Theorem 5.2.1], [14, §A.7]. We use below the facts that

- A is compactly supported iff there exist $\chi_1, \chi_2 \in \mathcal{D}$ such that $A = \chi_1 A \chi_2$, thus
- if $\chi_1, \chi_2 \in \mathcal{D}$ are compactly supported functions, then $\chi_1 A \chi_2$ is compactly supported, and
- if P is a differential operator and $\chi \in \mathcal{D}$, then both χP and $P \chi$ are compactly supported.

Ellipticity. We say that $B \in \Psi_{\hbar}^m$ is *elliptic* on $X \subset \overline{T}^* \mathbb{R}^d$ if there exists $c > 0$, independent of \hbar , such that

$$\langle \xi \rangle^{-m} |\sigma_{\hbar}(B)(x, \xi)| \geq c, \quad \text{for all } (x, \text{RC}(\xi)) \in X \text{ and for all } 0 < \hbar \leq \hbar_0. \quad (4.17)$$

A key feature of elliptic operators is that they are microlocally invertible; this is reflected in the following result.

Proposition 4.2 (Elliptic parametrix [14, Proposition E.32].) ² Let $A \in \Psi_{\hbar}^m$ and $B \in \Psi_{\hbar}^{\ell}$ be such that B is elliptic on $\text{WF}_{\hbar}(A)$. Then there exist $Q, Q' \in \Psi_{\hbar}^{m-\ell}$ such that

$$A = BQ + O(\hbar^{\infty})_{\Psi^{-\infty}} = Q'B + O(\hbar^{\infty})_{\Psi^{-\infty}}.$$

Theorem 4.3 (Elliptic estimate [14, Theorem E.33].) ² Let $A \in \Psi_{\hbar}^{m_1}$, $B_1 \in \Psi_{\hbar}^{m_2}$, and $P \in \Psi_{\hbar}^{\ell}$ be so that $B_1 P$ is elliptic on $\text{WF}_{\hbar}(A)$.

(i) Given $s, N > 0$, and $M > 0$, if $v \in \mathcal{D}'$ and $B_1 P v \in H^{s-m_2-\ell}$ then $Av \in H^{s-m_1}$ and there exists $C_s > 0$, $C_{N,M,s} > 0$ (independent of v and \hbar) such that

$$\|Av\|_{H_{\hbar}^{s-m_1}} \leq C_s \|B_1 P v\|_{H_{\hbar}^{s-m_2-\ell}} + C_{N,M,s} \hbar^M \|v\|_{H_{\hbar}^{-N}}. \quad (4.18)$$

(ii) If, in addition, A and $B_1 P$ are compactly supported, then there exists $\tilde{\chi} \in C_{\text{comp}}^{\infty}$ so that

$$\|Av\|_{H_{\hbar}^{s-m_1}} \leq C_s \|B_1 P v\|_{H_{\hbar}^{s-m_2-\ell}} + C_{N,M,s} \hbar^M \|\tilde{\chi} v\|_{H_{\hbar}^{-N}}. \quad (4.19)$$

Part (i) of Theorem 4.3 is proved by using Proposition 4.2 with $B = B_1 P \in \Psi_{\hbar}^{m_2+\ell}$, applying the resulting operator equation to v , and taking norms. The operator $Q' \in \Psi_{\hbar}^{m_1-m_2-\ell}$ and the constant C_s is then $\|Q'\|_{H_{\hbar}^{s-m_2-\ell} \rightarrow H_{\hbar}^{s-m_1}}$. The proof of Part (ii) is similar, using that, since A and $B_1 P$ are both compactly supported, there exists $\tilde{\chi} \in C_{\text{comp}}^{\infty}$ such that $(A - B_1 P)v = (A - B_1 P)\tilde{\chi}v$.

5 Proof of Theorem 3.1

In the notation introduced in §4, Theorem 3.1 becomes the following.

Theorem 5.1 Let A and n satisfy Assumption 2.1 and let $R > 0$ be such that $\text{supp}(1-A) \cup \text{supp}(1-n) \Subset B_R$. Given $f \in L^2(B_R)$, let u satisfy $P_{\hbar} u = \hbar^2 f$ in \mathbb{R}^d and the Sommerfeld radiation condition (1.2). Assume that, given $k_0 > 0$, $C_{\text{sol}}(k)$ is polynomially bounded (in the sense of Definition 2.6) for $k \in K \subset [k_0, \infty)$. Given $k_0 > 0$, let $\hbar_0 := k_0^{-1}$, and let $H := \{k^{-1} : k \in K\} \subset (0, \hbar_0]$.

Then there exist $C_3, C_4, C_5 > 0$ such that

$$u|_{B_R} = u_{H^2} + u_{\mathcal{A}}$$

where $u_{H^2} \in H_{\hbar}^2(B_R)$ with

$$\|\partial^{\alpha} u_{H^2}\|_{L^2(B_R)} \leq C_3 \hbar^{2-|\alpha|} \|f\|_{L^2(B_R)} \quad \text{for all } |\alpha| \leq 2 \text{ and for all } \hbar \in H \subset (0, \hbar_0], \quad (5.1)$$

and $u_{\mathcal{A}} \in C^{\infty}(B_R)$ with

$$\|\partial^{\beta} u_{\mathcal{A}}\|_{L^2(B_R)} \leq C_{\text{sol}}(\hbar^{-1}; R+2) C_4 \left(\frac{\hbar}{C_5}\right)^{1-|\beta|} \|f\|_{L^2(B_R)} \quad \text{for all } \beta \text{ and for all } \hbar \in H \subset (0, \hbar_0], \quad (5.2)$$

where C_3, C_4 , and C_5 depend on A, n, d , and \hbar_0 , but are independent of \hbar, f, α , and β .

²We highlight that working in \mathbb{R}^d (as opposed to on a general manifold defined by coordinate charts) allows us to remove the proper-support assumption appearing in [14, Proposition E.32, Theorem E.33].

5.1 Step 0: Restatement of bounds on the solution operator in semiclassical notation

The definition of C_{sol} (Definition 2.5) implies that, in semiclassical notation,

$$\|u\|_{H_{\hbar}^1(B_R)} \leq \hbar C_{\text{sol}}(\hbar^{-1}) \|f\|_{L^2(B_R)} \quad \text{for all } \hbar > 0. \quad (5.3)$$

It is convenient to record here in semiclassical notation the bound on the solution operator when C_{sol} is polynomially bounded.

Lemma 5.2 (Polynomial boundedness rewritten in terms of \hbar) *Given $f \in L_{\text{comp}}^2(\mathbb{R}^d)$, let $u \in H_{\text{loc}}^1(\mathbb{R}^d)$ be the solution to*

$$P_{\hbar}u = \hbar^2 f$$

satisfying the Sommerfeld radiation condition (1.2) (with $k = \hbar^{-1}$).

If $C_{\text{sol}}(k)$ is polynomially bounded for $k \in K \subset [k_0, \infty)$ (in the sense of Definition 2.6), then there exists $M > 0$ (independent of \hbar) such that, given $\chi \in C_{\text{comp}}^\infty(\mathbb{R}^d)$, there exists $C > 0$ (independent of \hbar but dependent on χ) such that

$$\|\chi u\|_{L^2} \leq C \hbar^{1-M} \|f\|_{L^2} \quad \text{for } \hbar \in H \subset (0, \hbar_0], \quad (5.4)$$

where $\hbar_0 := k_0^{-1}$ and $H := \{k^{-1} : k \in K\}$.

The bound (5.4) also holds with $\|\chi u\|_{L^2}$ replaced by $\|\chi u\|_{H_{\hbar}^1}$, but we only need it in the form (5.4) for what follows.

5.2 Step 1: The definitions of $u_{\mathcal{A}}$ and u_{H^2} .

The cut-off functions χ and χ_{μ} . Let $\chi \in C_{\text{comp}}^\infty(\mathbb{R}^d; [0, 1])$ be such that

$$\chi = \begin{cases} 1 & \text{in } B_1 \\ 0 & \text{outside } B_2. \end{cases} \quad (5.5)$$

For $\mu > 0$, let

$$\chi_{\mu}(\cdot) := \chi\left(\frac{\cdot}{\mu}\right). \quad (5.6)$$

We define $\mu_0 = \mu_0(\mathbf{A}, n)$ by

$$\mu_0(\mathbf{A}, n) := \left(1 + \frac{2n_{\text{max}}}{A_{\text{min}}}\right). \quad (5.7)$$

The reason for this definition is that it implies that

$$\text{if } |\xi|^2 \geq \mu_0 \quad \text{then} \quad \langle \xi \rangle^{-2} \sigma_{\hbar}(P) \geq \frac{A_{\text{min}}}{2} > 0. \quad (5.8)$$

Indeed, by (4.12),

$$\langle \xi \rangle^{-2} \sigma_{\hbar}(P) \geq \frac{A_{\text{min}} |\xi|^2 - n_{\text{max}}}{1 + |\xi|^2} = \frac{A_{\text{min}}}{2} + \left(\frac{A_{\text{min}}}{2}\right) \left(\frac{|\xi|^2 - 1 - 2n_{\text{max}}/A_{\text{min}}}{1 + |\xi|^2}\right),$$

and (5.8) follows. The importance of the property (5.8) is explained at the end of this subsection.

The frequency cut-offs Π_L and Π_H . We define Π_L and Π_H , the projections on low and high frequencies respectively, by (1.9) and (1.10). The definition of the quantisation Op_{\hbar} (4.8) and the change of variable $\zeta = \xi/\hbar$ imply that

$$\Pi_L = \text{Op}_{\hbar}(\chi_{\mu}(|\xi|^2)) \quad (5.9)$$

and

$$\Pi_H = I - \Pi_L. \quad (5.10)$$

These definitions and the definition of $\Psi_{\hbar}^m(\mathbb{R}^d)$ in §4 imply that $\Pi_L \in \Psi_{\hbar}^{-\infty}(\mathbb{R}^d)$ and $\Pi_H \in \Psi_{\hbar}^0(\mathbb{R}^d)$.

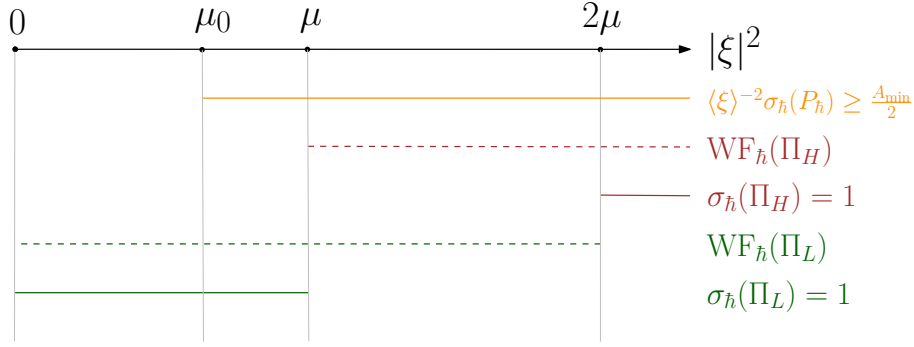


Figure 5.1: The locations of $\text{WF}_{\hbar}(\Pi_H)$ and $\text{WF}_{\hbar}(\Pi_L)$, the regions where the principal symbols of Π_H and Π_L equal one, and the region where P_{\hbar} is elliptic.

The locations of the wavefront sets of the frequency cut-offs, and the regions where their symbols equal one. In Figure 5.1 we show, as functions of $|\xi|^2$, the locations of $\text{WF}_{\hbar}(\Pi_H)$ and $\text{WF}_{\hbar}(\Pi_L)$, and the regions where $\sigma_{\hbar}(\Pi_H)$, and $\sigma_{\hbar}(\Pi_L)$ equal one. These locations/regions are obtained using (4.15) and (4.10) respectively. For example, since $1 - \chi_{\mu}(|\xi|^2) = 1$ for $|\xi|^2 \geq 2\mu$ and $= 0$ for $|\xi|^2 \leq \mu$, (4.10) and (4.15) imply that

$$\sigma_{\hbar}(\Pi_H) = 1 \text{ on } \{\xi : |\xi|^2 \geq 2\mu\} \quad \text{and} \quad \text{WF}_{\hbar}(\Pi_H) \subset \{\xi : |\xi|^2 \geq \mu\}. \quad (5.11)$$

We also record the following key consequence of the results summarised in Figure 5.1.

Lemma 5.3 *If $\mu \geq \mu_0$, then P_{\hbar} is elliptic on $\text{WF}_{\hbar}(\Pi_H)$.*

This property is central to our proof of the bound (5.1) on u_{H^2} , i.e., the high-frequency component. It is a consequence of (5.8), and the reason why we choose μ_0 as in (5.7) is for this ellipticity result to hold.

The definitions of $u_{\mathcal{A}}$ and u_{H^2} . As described in §1.2, we choose $\varphi \in C_{\text{comp}}^{\infty}(\mathbb{R}^d)$ be equal to one on B_{R+1} and vanish outside B_{R+2} . We then let

$$w := \varphi u$$

and we define

$$u_{\mathcal{A}} := (\Pi_L w)|_{B_R} \quad \text{and} \quad u_{H^2} := (\Pi_H w)|_{B_R}.$$

5.3 Step 2: Proof of the bound (5.2) on $u_{\mathcal{A}}$ (the low-frequency component)

Since $\Pi_L \in \Psi_{\hbar}^{-\infty}$, Part (iii) of Theorem 4.1, together with Sobolev embedding, gives $\Pi_L w \in C^{\infty}$.

The definition of Π_L (1.9) and Plancherel's identity (4.4) for the standard (i.e. non semiclassical) Fourier transform imply that

$$\|\partial^{\beta}(\Pi_L w)\|_{L^2} = \frac{1}{(2\pi)^{d/2}} \|(\cdot)^{\beta} \mathcal{F}(\Pi_L w)(\cdot)\|_{L^2} = \frac{1}{(2\pi)^{d/2}} \|(\cdot)^{\beta} \chi_{\mu}(\hbar^2 |\cdot|^2) \mathcal{F}w(\cdot)\|_{L^2}. \quad (5.12)$$

The definitions of χ (5.5) and χ_{μ} (5.6) imply that $\chi_{\mu}(\xi) = 0$ for $|\xi| \geq 2\mu$, so

$$\chi_{\mu}(\hbar^2 |\zeta|^2) = 0 \quad \text{for } |\zeta| \geq \sqrt{2\mu} \hbar^{-1}.$$

Using this fact, and then (in this order) the fact that $|\chi_{\mu}| \leq 1$, Plancherel's identity for the standard Fourier transform, the fact that $\varphi = 0$ outside B_{R+2} , and the definition of C_{sol} (2.12), we find from (5.12) that

$$\|\partial^{\beta}(\Pi_L \varphi u)\|_{L^2} \leq \frac{(2\mu)^{|\beta|/2}}{(2\pi)^{d/2}} \hbar^{-|\beta|} \|\chi_{\mu}(\hbar^2 |\cdot|^2) \mathcal{F}(\varphi u)(\cdot)\|_{L^2}$$

$$\begin{aligned}
&\leq \frac{(2\mu)^{|\beta|/2}}{(2\pi)^{d/2}} \hbar^{-|\beta|} \|\mathcal{F}(\varphi u)\|_{L^2} \\
&\leq (2\mu)^{|\beta|/2} \hbar^{-|\beta|} \|\varphi u\|_{L^2} \\
&\leq (2\mu)^{|\beta|/2} \hbar^{-|\beta|} \hbar C_{\text{sol}}(\hbar^{-1}; R+2) \|f\|_{L^2(B_R)}.
\end{aligned}$$

Since

$$\|\partial^\beta u_{\mathcal{A}}\|_{L^2(B_R)} = \|\partial^\beta(\Pi_L w)\|_{L^2(B_R)} \leq \|\partial^\beta(\Pi_L w)\|_{L^2},$$

the bound (5.2) then follows with $C_4 := \sqrt{2\mu}$ and $C_5 := \sqrt{2\mu}$.

5.4 Step 3: Proof of the bound (5.1) on u_{H^2} (the high-frequency component)

By the inequality (4.6), it is sufficient to prove that

$$\|\Pi_H w\|_{H_{\hbar}^2} \leq C_3 \hbar^2 \|f\|_{L^2(B_R)} \quad \text{for all } \hbar \in H \subset (0, \hbar_0). \quad (5.13)$$

It is instructive to first prove (5.13) under the assumption that $C_{\text{sol}}(k) \lesssim 1$ (which, by Theorem 2.7 is ensured if \mathbf{A} and n are nontrapping). Indeed, as discussed in §1.2, this proof only requires that P_{\hbar} is elliptic on $\text{WF}_{\hbar}(\Pi_H)$; i.e., Lemma 5.3. Throughout the rest of this section, therefore, we assume that $\mu \geq \mu_0$, so that the result of Lemma 5.3 holds.

5.4.1 Proof of (5.13) under the assumption that $C_{\text{sol}}(k) \lesssim 1$

We seek to apply Part (i) of Theorem 4.3 with $A = \Pi_H$ (so $m_1 = 0$), $B_1 = 1$ (so $m_2 = 0$), and $P = P_{\hbar}$ (so $\ell = 2$). By Lemma 5.3, $B_1 P$ is elliptic on $\text{WF}_{\hbar}(A)$. We can therefore apply Theorem 4.3 and obtain that, given $N, N' > 0$,

$$\|\Pi_H w\|_{H_{\hbar}^2} \lesssim \|P_{\hbar} w\|_{L^2} + \hbar^{N'} \|w\|_{H_{\hbar}^{-N}}, \quad (5.14)$$

where the omitted constant in \lesssim depends on N and N' . Since $P_{\hbar} u = \hbar^2 f$,

$$P_{\hbar} w = [P_{\hbar}, \varphi]u + \hbar^2 \varphi f,$$

where $[\cdot, \cdot]$ is the standard commutator defined by $[A_1, A_2] := A_1 A_2 - A_2 A_1$, so that (5.14) becomes

$$\|\Pi_H w\|_{H_{\hbar}^2} \lesssim \|[P_{\hbar}, \varphi]u\|_{L^2} + \hbar^2 \|f\|_{L^2} + \hbar^{N'} \|w\|_{H_{\hbar}^{-N}}. \quad (5.15)$$

Direct calculation, using the fact that $\text{supp } \varphi \subset B_{R+2}$, implies that

$$\|[P_{\hbar}, \varphi]u\|_{L^2} \lesssim \hbar \|u\|_{H_{\hbar}^1(B_{R+2})}, \quad (5.16)$$

where the omitted constant depends on φ , and hence on R .

Combining (5.15) and (5.16), and recalling that $\text{supp } \varphi \subset B_{R+2}$, we have

$$\|\Pi_H w\|_{H_{\hbar}^2} \lesssim \hbar \|u\|_{H_{\hbar}^1(B_{R+2})} + \hbar^2 \|f\|_{L^2(B_R)} + \hbar^{N'} \|u\|_{H_{\hbar}^{-N}(B_{R+2})}.$$

Choosing $N = 0$ and $N' = 1$, and then using (5.3), we obtain

$$\|\Pi_H w\|_{H_{\hbar}^2} \lesssim \hbar^2 \left(1 + C_{\text{sol}}(\hbar^{-1})\right) \|f\|_{L^2(B_R)}. \quad (5.17)$$

If $C_{\text{sol}}(\hbar^{-1}) \lesssim 1$, then this implies (5.13). However, if $C_{\text{sol}}(\hbar^{-1}) \gg 1$ (as occurs when C_{sol} is polynomially bounded in the sense of Definition 2.6 with $M > 0$) then (5.17) is a weaker bound than (5.13).

5.4.2 Proof of (5.13) under the assumption that $C_{\text{sol}}(k)$ is polynomially bounded

Inspecting the argument in §5.4.1, we see that the assumption that $C_{\text{sol}}(k) \lesssim 1$ is needed to get a good bound on the commutator term $[P_{\hbar}, \varphi]u$. To remove this commutator term, one idea is to use the elliptic estimate in Part (i) of Theorem 4.3, using the fact that P_{\hbar} is elliptic on $\text{WF}_{\hbar}(\Pi_H \varphi)$, and apply the estimate with $v := u$. However, the error term would not be compactly supported and we would be unable to control it using the polynomial bound on the solution operator (5.4). We therefore introduce additional spatial cut-offs on the left of $\Pi_H \varphi$ and P_{\hbar} to create compactly-supported operators and have a compactly-supported error term thanks to Part (ii) of Theorem 4.3.

To this end, let $\varphi_1, \varphi_2 \in C_{\text{comp}}^{\infty}(\mathbb{R}^d)$ be such that $\varphi_1 = 1$ on $\text{supp } \varphi$ and $\varphi_2 = 1$ on $\text{supp } \varphi_1$; we then write

$$\Pi_H \varphi u = (1 - \varphi_1) \Pi_H \varphi u + \varphi_1 \Pi_H \varphi u. \quad (5.18)$$

Since $1 - \varphi_1 = 0$ on $\text{supp } \varphi$, using (4.14) and (4.15), we obtain that

$$\text{WF}_{\hbar}((1 - \varphi_1) \Pi_H \varphi) \subset \bar{T}^*(\text{supp}(1 - \varphi_1)) \cap \bar{T}^*(\text{supp } \varphi) = \emptyset.$$

Hence, by (4.16), $(1 - \varphi_1) \Pi_H \varphi = O(\hbar^{\infty})_{\Psi^{-\infty}}$, and, by the definition of the residual class (4.9), for any $N \geq 1$ there exists $C_N > 0$ so that

$$\|(1 - \varphi_1) \Pi_H \varphi u\|_{H_{\hbar}^2} = \|(1 - \varphi_1) \Pi_H \varphi \varphi_1 u\|_{H_{\hbar}^2} \leq C_N \hbar^N \|\varphi_1 u\|_{L^2}, \quad (5.19)$$

where we used the fact that $\varphi_1 = 1$ on $\text{supp } \varphi$ in the first equality.

It therefore remains to control $\varphi_1 \Pi_H \varphi u$; to do this, we use the elliptic estimate of Theorem 4.3.

Lemma 5.4 $\varphi_2 P_{\hbar}$ is elliptic on $\text{WF}_{\hbar}(\varphi_1 \Pi_H \varphi)$.

Proof. By (4.14) and (4.15), $\text{WF}_{\hbar}(\varphi_1 \Pi_H \varphi) \subset \bar{T}^*(\text{supp } \varphi_1) \cap \text{WF}_{\hbar} \Pi_H$. Since $\varphi_2 = 1$ on $\text{supp } \varphi_1$, the result is a direct consequence of Lemma 5.3. \blacksquare

By the facts about compactly-supported operators recalled in §4, $\varphi_1 \Pi_H \varphi$ and $\varphi_2 P_{\hbar}$ are compactly supported. Therefore, by Lemma 5.4, we can apply Part (ii) of Theorem 4.3 with $A = \varphi_1 \Pi_H \varphi$, $B_1 = \varphi_2$, $P = P_{\hbar}$, $m_1 = 0$, $m_2 = 0$, $\ell = 2$. This result implies that there exists $\tilde{\chi} \in C_{\text{comp}}^{\infty}$, and, for any $N' \geq 1$, there exists $C_{N'} > 0$ such that

$$\|\varphi_1 \Pi_H \varphi u\|_{H_{\hbar}^2} \lesssim \|\varphi_2 P_{\hbar} u\|_{L^2} + C_{N'} \hbar^{N'} \|\tilde{\chi} u\|_{L^2} = \hbar^2 \|\varphi_2 f\|_{L^2} + C_{N'} \hbar^{N'} \|\tilde{\chi} u\|_{L^2}. \quad (5.20)$$

Collecting (5.18), (5.19), (5.20), using (5.4), and choosing $N = N' = M + 1$, we obtain (5.13).

6 Proof of Theorem 3.4

The two ingredients for the proof of Theorem 3.4 are

- Lemma 6.4, which is the standard duality argument giving a condition for quasi-optimality to hold in terms of how well the solution of the adjoint problem is approximated by the finite-element space (measured by the quantity $\eta(V_N)$ defined by (6.3)), and
- Lemma 6.5 that bounds $\eta(V_N)$ using the decomposition from Theorem 3.1.

Regarding Lemma 6.4: we recall that this argument came out of ideas introduced in [43], was then formalised in [42], and has been used extensively in the analysis of the Helmholtz FEM; see, e.g., [1, 26, 33, 42, 35, 36, 51, 50, 12, 9, 31, 10, 17, 21, 16].

Before stating Lemma 6.4 we need to introduce some notation.

Definition 6.1 (The adjoint sesquilinear form $a^*(\cdot, \cdot)$) The adjoint sesquilinear form, $a^*(u, v)$, to the sesquilinear form $a(\cdot, \cdot)$ defined in (2.4) is given by

$$a^*(u, v) := \overline{a(v, u)} = \int_{B_R} \left((A \nabla u) \cdot \overline{\nabla v} - k^2 n u \bar{v} \right) - \langle \gamma u, \text{DtN}_k(\gamma v) \rangle_{\partial B_R}.$$

A key role is played by the solution operator of the adjoint variational problem with data in $L^2(B_R)$; we therefore introduce the following notation.

Definition 6.2 (Adjoint solution operator \mathcal{S}^*) Given $f \in L^2(B_R)$, let \mathcal{S}^*f be defined as the solution of the variational problem

$$\text{find } \mathcal{S}^*f \in H^1(B_R) \quad \text{such that} \quad a^*(\mathcal{S}^*f, v) = \int_{B_R} f \bar{v} \quad \text{for all } v \in H^1(B_R). \quad (6.1)$$

Green's second identity applied to solutions of the Helmholtz equation satisfying the Sommerfeld radiation condition (1.2) implies that $\langle \text{DtN}_k \psi, \bar{\phi} \rangle_{\partial B_R} = \langle \text{DtN}_k \phi, \bar{\psi} \rangle_{\partial B_R}$ (see, e.g., [46, Lemma 6.13]); thus $a(\bar{v}, u) = a(\bar{u}, v)$ and so the definition (6.1) implies that

$$a(\overline{\mathcal{S}^*f}, v) = (\bar{f}, v)_{L^2(B_R)} \quad \text{for all } v \in H^1(B_R). \quad (6.2)$$

Definition 6.3 ($\eta(V_N)$) Given a sequence of finite-dimensional spaces $(V_N)_{N=0}^\infty$ (as described in §2.3), let

$$\eta(V_N) := \sup_{0 \neq f \in L^2(B_R)} \min_{v_N \in V_N} \frac{\|\mathcal{S}^*f - v_N\|_{H_k^1(B_R)}}{\|f\|_{L^2(B_R)}}. \quad (6.3)$$

Lemma 6.4 (Conditions for quasi-optimality) If

$$k \eta(V_N) \leq \frac{1}{C_{\text{cont}}} \sqrt{\frac{A_{\min}}{2(n_{\max} + A_{\min})}},$$

then the Galerkin equations (2.15) have a unique solution which satisfies

$$\|u - u_h\|_{H_k^1(B_R)} \leq \frac{2C_{\text{cont}}}{A_{\min}} \left(\min_{v_N \in V_N} \|u - v_N\|_{H_k^1(B_R)} \right).$$

Proof. Using the inequality (2.10), we see that $a(\cdot, \cdot)$ satisfies the Gårding inequality

$$\Re(a(v, v)) \geq A_{\min} \|v\|_{H_k^1(B_R)}^2 - 2k^2(n_{\max} + A_{\min}) \|v\|_{L^2(B_R)}^2 \quad (6.4)$$

and the result follows from, e.g., the account [46, Theorem 6.32] of the standard duality argument with (in the notation of [46]) $\alpha = A_{\min}$ and $C_{\mathcal{V}} = 2k^2(n_{\max} + A_{\min})$. ■

Lemma 6.5 (Bound on $\eta(V_N)$ using the decomposition from Theorem 3.1) Let A and n satisfy Assumption 2.1 and let $R > 0$ be such that $\text{supp}(1 - A) \cup \text{supp}(1 - n) \Subset B_R$. Let $(V_N)_{N=0}^\infty$ be the piecewise-polynomial approximation spaces described in [35, §5]. There exists $C_6, C_7, \sigma > 0$, all independent of k, h , and p , such that

$$k \eta(V_N) \leq C_6 C_3 \frac{hk}{p} \left(1 + \frac{kh}{p} \right) + C_7 C_{\text{sol}}(k) \left[\left(\frac{h}{h + \sigma} \right)^p \left(1 + \frac{hk}{h + \sigma} \right) + k \left(\frac{kh}{\sigma p} \right)^p \left(\frac{1}{p} + \frac{kh}{\sigma p} \right) \right]. \quad (6.5)$$

The constants C_6 and σ only depend on the constants in [35, Assumption 5.2] defining the element maps from the reference element; C_7 depends on these constants, and additionally on C_5 .

Proof. This proof is very similar to the proof of [35, Theorem 5.5]. Indeed, [35, Theorem 5.5] proves a bound very similar to (6.5) starting from bounds almost identical to the bounds (3.1) and (3.2) (recalling Remark 3.3 about notation). The only difference is that the bound (3.2) contains C_{sol} , which depends on k (whereas in [35] $C_{\text{sol}} \sim 1$), and so we now need to keep track of how C_{sol} enters the proof of [35, Theorem 5.5].

From the definition (6.3), it is sufficient to show that, given $f \in L^2(B_R)$, there exists $w_N \in V_N$ such that

$$\|\mathcal{S}^*f - w_N\|_{H_k^1(B_R)} \leq C \|f\|_{L^2(B_R)}, \quad (6.6)$$

where C is the right-hand side of (6.5) divided by k . Let $v := \mathcal{S}^*f$; by (6.2) and Part (i) of Lemma 2.3, \bar{v} satisfies the assumptions of Theorem 3.1 with f replaced by \bar{f} , and so the bounds (3.1) and (3.2) hold with u replaced by v .

By [35, First equation on Page 1896] (which uses [35, Theorem B.4]), the bound (3.6) holds, and thus there exists $w_N^{(1)} \in V_N$ such that

$$\|v_{H^2} - w_N^{(1)}\|_{H_k^1(B_R)} \leq C_6 \frac{h}{p} \left(1 + \frac{kh}{p}\right) |v|_{H^2(B_R)}$$

and so

$$\|v_{H^2} - w_N^{(1)}\|_{H_k^1(B_R)} \leq C_6 \frac{h}{p} \left(1 + \frac{kh}{p}\right) C_3 \|f\|_{L^2(B_R)} \quad (6.7)$$

by (3.1).

For the approximation of v_A , the only change to the argument in [35] is that a multiplicative factor of $(C_{\text{sol}})^2$ must be included on the right-hand side of [35, Equation 5.8]. Then [35, Equations 5.8 and 5.9] implies that there exists C_7 and $w_N^{(2)} \in V_N$ such that

$$k \|v_A - w_N^{(2)}\|_{H_k^1(B_R)} \leq C_7 C_{\text{sol}}(k) \left[\left(\frac{h}{h+\sigma}\right)^p \left(1 + \frac{hk}{h+\sigma}\right) + k \left(\frac{kh}{\sigma p}\right)^p \left(\frac{1}{p} + \frac{kh}{\sigma p}\right) \right] \|f\|_{L^2(B_R)} \quad (6.8)$$

(observe that this equation is identical to [35, Last equation on Page 1896] except for the factor C_{sol} on the right-hand side).

Let $w_N := w_N^{(1)} + w_N^{(2)}$. By the triangle inequality, the decomposition $v = v_{H^2} + v_A$ on B_R , and the inequalities (6.7) and (6.8), the inequality (6.6) holds with C the right-hand side of (6.5) and the proof is complete. \blacksquare

Corollary 6.6 (Conditions under which $k\eta(V_N)$ is arbitrarily small) *Let the assumptions of Lemma 6.5 hold. Given $\varepsilon > 0$ and $k_0 > 0$, there exists $\mathcal{C}_1, \mathcal{C}_2 > 0$, depending only on $\varepsilon, \mathcal{C}_3, \mathcal{C}_6, \mathcal{C}_7, \sigma$, and k_0 , such that if*

$$\frac{hk}{p} \leq \mathcal{C}_1 \quad \text{and} \quad p \geq \mathcal{C}_2 \left(1 + \log k + \log(C_{\text{sol}}(k))\right),$$

then

$$k\eta(V_N) \leq \varepsilon \quad \text{for all } k \geq k_0.$$

Proof. This proof is essentially identical to the proofs of [35, Corollary 5.6] and [36, Theorem 5.8]. First choose \mathcal{C}_1 sufficiently small such that $\mathcal{C}_1 < \sigma$ and

$$C_6 C_3 \mathcal{C}_1 (1 + \mathcal{C}_1) \leq \frac{\varepsilon}{2}$$

From the bound on $k\eta(V_N)$ (6.5), it is then sufficient to show that

$$C_7 C_{\text{sol}}(k) \left[\left(\frac{h}{h+\sigma}\right)^p \left(1 + \frac{hk}{h+\sigma}\right) + k \left(\frac{kh}{\sigma p}\right)^p \left(\frac{1}{p} + \frac{kh}{\sigma p}\right) \right] \quad (6.9)$$

can be made $\leq \varepsilon/2$. Let

$$\theta_1 := \frac{h}{h+\sigma} \quad \text{and} \quad \theta_2 := \frac{\mathcal{C}_1}{\sigma},$$

so that (6.9) is bounded by

$$C_7 C_{\text{sol}}(k) \left[(\theta_1)^p \left(1 + \frac{\mathcal{C}_1 p}{\sigma}\right) + k (\theta_2)^p \left(\frac{1}{p} + \frac{\mathcal{C}_1}{\sigma}\right) \right];$$

the result then follows since $\theta_1, \theta_2 < 1$. \blacksquare

Proof of Theorem 3.4. This follows by combining Lemma 6.4 and Corollary 6.6. \blacksquare

Acknowledgements

The authors thank Martin Averseng (ETH Zürich) and an anonymous referee for highlighting simplifications of the arguments in a earlier version of the paper. We also thank Théophile Chaumont-Frelet (INRIA, Nice) for useful discussions about the results of [35], [36]. DL and EAS acknowledge support from EPSRC grant EP/1025995/1. JW was partly supported by Simons Foundation grant 631302.

References

- [1] A. K. AZIZ, R. B. KELLOGG, AND A. B. STEPHENS, *A two point boundary value problem with a rapidly oscillating solution*, Numer. Math., 53 (1988), pp. 107–121.
- [2] I. M. BABUŠKA AND S. A. SAUTER, *Is the pollution effect of the FEM avoidable for the Helmholtz equation considering high wave numbers?*, SIAM Review, (2000), pp. 451–484.
- [3] H. BARUCQ, T. CHAUMONT-FRELET, AND C. GOUT, *Stability analysis of heterogeneous Helmholtz problems and finite element solution based on propagation media approximation*, Math. Comp., 86 (2017), pp. 2129–2157.
- [4] M. BERNKOPF, T. CHAUMONT-FRELET, AND J. M. MELENK, *Stability and convergence of Galerkin discretizations of the Helmholtz equation in piecewise smooth media*, https://numericalwaves.sciencesconf.org/data/program/abstract_melenk.pdf, (2020).
- [5] N. BURQ, *Décroissance des ondes absence de de l'énergie locale de l'équation pour le problème extérieur et absence de résonance au voisinage du réel*, Acta Math., 180 (1998), pp. 1–29.
- [6] ———, *Semi-classical estimates for the resolvent in nontrapping geometries*, International Mathematics Research Notices, 2002 (2002), pp. 221–241.
- [7] S. N. CHANDLER-WILDE AND P. MONK, *Wave-number-explicit bounds in time-harmonic scattering*, SIAM J. Math. Anal., 39 (2008), pp. 1428–1455.
- [8] T. CHAUMONT-FRELET, *On high order methods for the heterogeneous Helmholtz equation*, Computers & Mathematics with Applications, 72 (2016), pp. 2203–2225.
- [9] T. CHAUMONT-FRELET AND S. NICAISE, *High-frequency behaviour of corner singularities in Helmholtz problems*, ESAIM: Math. Model. Numer. Anal., 52 (2018), pp. 1803–1845.
- [10] ———, *Wavenumber explicit convergence analysis for finite element discretizations of general wave propagation problem*, IMA J. Numer. Anal., 40 (2020), pp. 1503–1543.
- [11] M. COSTABEL, M. DAUGE, AND S. NICAISE, *Corner Singularities and Analytic Regularity for Linear Elliptic Systems. Part I: Smooth domains.*, (2010). https://hal.archives-ouvertes.fr/file/index/docid/453934/filename/CoDaNi_Analytic_Part_I.pdf.
- [12] Y. DU AND H. WU, *Preasymptotic error analysis of higher order FEM and CIP-FEM for Helmholtz equation with high wave number*, SIAM J. Numer. Anal., 53 (2015), pp. 782–804.
- [13] J. J. DUISTERMAAT AND L. HÖRMANDER, *Fourier integral operators. ii*, Acta mathematica, 128 (1972), pp. 183–269.
- [14] S. DYATLOV AND M. ZWORSKI, *Mathematical theory of scattering resonances*, vol. 200 of Graduate Studies in Mathematics, American Mathematical Society, 2019.
- [15] S. ESTERHAZY AND J. M. MELENK, *On stability of discretizations of the Helmholtz equation*, in Numerical Analysis of Multiscale Problems, I. G. Graham, T. Y. Hou, O. Lakkis, and R. Scheichl, eds., Springer, 2012, pp. 285–324.
- [16] J. GALKOWSKI, E. A. SPENCE, AND J. WUNSCH, *Optimal constants in nontrapping resolvent estimates*, Pure and Applied Analysis, 2 (2020), pp. 157–202.
- [17] D. GALLISTL, T. CHAUMONT-FRELET, S. NICAISE, AND J. TOMEZYK, *Wavenumber explicit convergence analysis for finite element discretizations of time-harmonic wave propagation problems with perfectly matched layers*, hal preprint 01887267, (2018).
- [18] M. GANESH AND C. MORGENSTERN, *A coercive heterogeneous media Helmholtz model: formulation, wavenumber-explicit analysis, and preconditioned high-order FEM*, Numerical Algorithms, (2019), pp. 1–47.
- [19] S. GONG, I. G. GRAHAM, AND E. A. SPENCE, *Domain decomposition preconditioners for high-order discretizations of the heterogeneous Helmholtz equation*, IMA J. Num. Anal., 41 (2021), pp. 2139–2185.
- [20] I. G. GRAHAM, O. R. PEMBERY, AND E. A. SPENCE, *The Helmholtz equation in heterogeneous media: a priori bounds, well-posedness, and resonances*, Journal of Differential Equations, 266 (2019), pp. 2869–2923.
- [21] I. G. GRAHAM AND S. A. SAUTER, *Stability and finite element error analysis for the Helmholtz equation with variable coefficients*, Math. Comp., 89 (2020), pp. 105–138.
- [22] B. HELFFER AND D. ROBERT, *Calcul fonctionnel par la transformation de Mellin et opérateurs admissibles*, J. Funct. Anal., 53 (1983), pp. 246–268.

- [23] B. HELFFER AND J. SJÖSTRAND, *Équation de Schrödinger avec champ magnétique et équation de Harper*, in Schrödinger operators (Sønderborg, 1988), vol. 345 of Lecture Notes in Phys., Springer, Berlin, 1989, pp. 118–197.
- [24] L. HÖRMANDER, *The Analysis of Linear Differential Operators. I, Distribution Theory and Fourier Analysis*, Springer-Verlag, Berlin, 1983.
- [25] G. C. HSIAO AND W. L. WENDLAND, *Boundary integral equations*, vol. 164 of Applied Mathematical Sciences, Springer, 2008.
- [26] F. IHLENBURG AND I. BABUŠKA, *Finite element solution of the Helmholtz equation with high wave number Part I: The h -version of the FEM*, Comput. Math. Appl., 30 (1995), pp. 9–37.
- [27] D. LAFONTAINE, E. A. SPENCE, AND J. WUNSCH, *Decompositions of high-frequency Helmholtz solutions via functional calculus, and application to the finite element method*, arXiv preprint arXiv:2102.13081, (2021).
- [28] ———, *For most frequencies, strong trapping has a weak effect in frequency-domain scattering*, Communications on Pure and Applied Mathematics, 74 (2021), pp. 2025–2063.
- [29] ———, *A sharp relative-error bound for the Helmholtz h -FEM at high frequency*, Numerische Mathematik, 150 (2022), pp. 137–178.
- [30] P. D. LAX AND R. S. PHILLIPS, *Scattering Theory*, Academic Press, revised ed., 1989.
- [31] Y. LI AND H. WU, *FEM and CIP-FEM for Helmholtz Equation with High Wave Number and Perfectly Matched Layer Truncation*, SIAM J. Numer. Anal., 57 (2019), pp. 96–126.
- [32] A. MARTINEZ, *An introduction to semiclassical and microlocal analysis*, vol. 994, Springer, 2002.
- [33] J. M. MELENK, *On generalized finite element methods*, PhD thesis, The University of Maryland, 1995.
- [34] J. M. MELENK, A. PARSANIA, AND S. SAUTER, *General DG-methods for highly indefinite Helmholtz problems*, Journal of Scientific Computing, 57 (2013), pp. 536–581.
- [35] J. M. MELENK AND S. SAUTER, *Convergence analysis for finite element discretizations of the Helmholtz equation with Dirichlet-to-Neumann boundary conditions*, Math. Comp, 79 (2010), pp. 1871–1914.
- [36] ———, *Wavenumber explicit convergence analysis for Galerkin discretizations of the Helmholtz equation*, SIAM J. Numer. Anal., 49 (2011), pp. 1210–1243.
- [37] J. C. NÉDÉLEC, *Acoustic and electromagnetic equations: integral representations for harmonic problems*, Springer Verlag, 2001.
- [38] O. R. PEMBERY, *The Helmholtz Equation in Heterogeneous and Random Media: Analysis and Numerics*, PhD thesis, University of Bath, 2020.
- [39] J. V. RALSTON, *Trapped rays in spherically symmetric media and poles of the scattering matrix*, Communications on Pure and Applied Mathematics, 24 (1971), pp. 571–582.
- [40] D. ROBERT, *Autour de l'approximation semi-classique*, vol. 68 of Progress in Mathematics, Birkhäuser Boston, Inc., Boston, MA, 1987.
- [41] J. SARANEN AND G. VAINIKKO, *Periodic integral and pseudodifferential equations with numerical approximation*, Springer, 2002.
- [42] S. A. SAUTER, *A refined finite element convergence theory for highly indefinite Helmholtz problems*, Computing, 78 (2006), pp. 101–115.
- [43] A. H. SCHATZ, *An observation concerning Ritz-Galerkin methods with indefinite bilinear forms*, Math. Comp., 28 (1974), pp. 959–962.
- [44] J. SJÖSTRAND, *A trace formula and review of some estimates for resonances*, in Microlocal analysis and spectral theory (Lucca, 1996), vol. 490 of NATO Adv. Sci. Inst. Ser. C Math. Phys. Sci., Kluwer Acad. Publ., Dordrecht, 1997, pp. 377–437.
- [45] J. SJÖSTRAND AND M. ZWORSKI, *Complex scaling and the distribution of scattering poles*, J. Amer. Math. Soc., 4 (1991), pp. 729–769.
- [46] E. A. SPENCE, *Overview of Variational Formulations for Linear Elliptic PDEs*, in Unified transform method for boundary value problems: applications and advances, A. S. Fokas and B. Pelloni, eds., SIAM, 2015, pp. 93–159.
- [47] M. E. TAYLOR, *Partial differential equations II, Qualitative studies of linear equations, volume 116 of Applied Mathematical Sciences*, Springer-Verlag, New York, 1996.
- [48] B. R. VAINBERG, *On the short wave asymptotic behaviour of solutions of stationary problems and the asymptotic behaviour as $t \rightarrow \infty$ of solutions of non-stationary problems*, Russian Mathematical Surveys, 30 (1975), pp. 1–58.
- [49] ———, *Asymptotic methods in equations of mathematical physics*, Gordon & Breach Science Publishers, New York, 1989. Translated from the Russian by E. Primrose.
- [50] H. WU, *Pre-asymptotic error analysis of CIP-FEM and FEM for the Helmholtz equation with high wave number. Part I: linear version*, IMA J. Numer. Anal., 34 (2014), pp. 1266–1288.
- [51] L. ZHU AND H. WU, *Preasymptotic error analysis of CIP-FEM and FEM for Helmholtz equation with high wave number. Part II: hp version*, SIAM J. Numer. Anal., 51 (2013), pp. 1828–1852.
- [52] M. ZWORSKI, *Semiclassical analysis*, vol. 138 of Graduate Studies in Mathematics, American Mathematical Society, Providence, RI, 2012.