

Ensuring quasi-optimality for the Helmholtz problem

Tim van Beeck¹, Umberto Zerbinati²

¹Institute for Numerical and Applied Mathematics,
University of Göttingen

²Mathematical Institute, University of Oxford

European Finite Element Fair 2024, London

Let Ω be Lipschitz. For a wave-number k , find u s.t.

$$\begin{aligned} -\Delta u - k^2 u &= f && \text{in } \Omega, \\ u &= 0 && \text{on } \partial\Omega. \end{aligned}$$

Pollution effect¹: For **fixed** mesh size h , we loose quasi-optimality as the wave-number k **increases**.

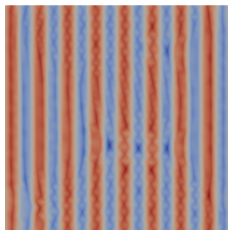
¹extensively studied, e.g. [Babuška, Sauter, 2000], [Melenk, Sauter, 2010 & 2011], [Bernkopf, Chaumont-Frelet, Melenk 2024]

Let Ω be Lipschitz. For a wave-number k , find u s.t.

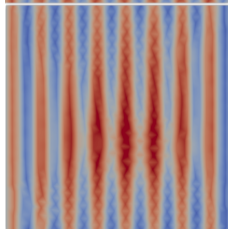
$$\begin{aligned} -\Delta u - k^2 u &= f && \text{in } \Omega, \\ u &= 0 && \text{on } \partial\Omega. \end{aligned}$$

Pollution effect¹: For **fixed** mesh size h , we lose quasi-optimality as the wave-number k **increases**.

$$u_{\text{ex}} = \sin(16\pi x), \quad k = 1$$



$$\Pi_h^{L^2} u_{\text{ex}} \\ h = 0.1$$



$$u_h \\ h = 0.1$$

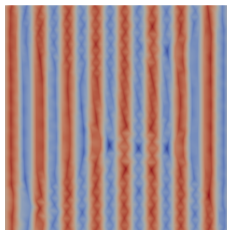
¹extensively studied, e.g. [Babuška, Sauter, 2000], [Melenk, Sauter, 2010 & 2011], [Bernkopf, Chaumont-Frelet, Melenk 2024]

Let Ω be Lipschitz. For a wave-number k , find u s.t.

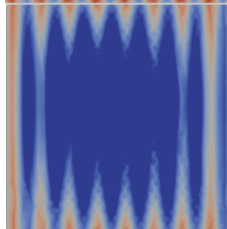
$$\begin{aligned} -\Delta u - k^2 u &= f && \text{in } \Omega, \\ u &= 0 && \text{on } \partial\Omega. \end{aligned}$$

Pollution effect¹: For **fixed** mesh size h , we lose quasi-optimality as the wave-number k **increases**.

$$u_{\text{ex}} = \sin(16\pi x), \quad k = 5$$



$$\Pi_h^{L^2} u_{\text{ex}} \\ h = 0.1$$



$$u_h \\ h = 0.1$$

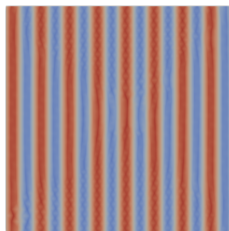
¹extensively studied, e.g. [Babuška, Sauter, 2000], [Melenk, Sauter, 2010 & 2011], [Bernkopf, Chaumont-Frelet, Melenk 2024]

Let Ω be Lipschitz. For a wave-number k , find u s.t.

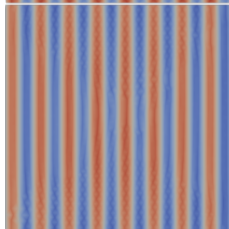
$$\begin{aligned} -\Delta u - k^2 u &= f && \text{in } \Omega, \\ u &= 0 && \text{on } \partial\Omega. \end{aligned}$$

Pollution effect¹: For **fixed** mesh size h , we lose quasi-optimality as the wave-number k **increases**.

$$u_{\text{ex}} = \sin(16\pi x), \quad k = 5$$



$$\Pi_h^{L^2} u_{\text{ex}} \\ h = 0.075$$



$$u_h \\ h = 0.075$$

¹extensively studied, e.g. [Babuška, Sauter, 2000], [Melenk, Sauter, 2010 & 2011], [Bernkopf, Chaumont-Frelet, Melenk 2024]

Let Ω be Lipschitz. For a wave-number k , find u s.t.

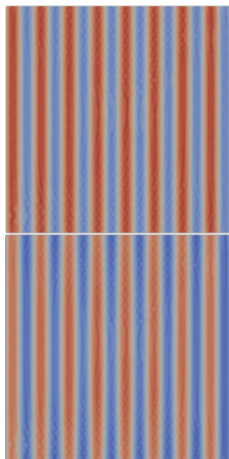
$$\begin{aligned} -\Delta u - k^2 u &= f && \text{in } \Omega, \\ u &= 0 && \text{on } \partial\Omega. \end{aligned}$$

Pollution effect¹: For **fixed** mesh size h , we loose quasi-optimality as the wave-number k **increases**.

Goal: Derive **practical** scheme to generate a mesh that guarantees quasi-optimality

¹extensively studied, e.g. [Babuška, Sauter, 2000], [Melenk, Sauter, 2010 & 2011], [Bernkopf, Chaumont-Frelet, Melenk 2024]

$$u_{\text{ex}} = \sin(16\pi x), \quad k = 5$$



$$\Pi_h^{L^2} u_{\text{ex}} \\ h = 0.075$$

$$u_h \\ h = 0.075$$

Let Ω be **Lipschitz**. For a wave-number k , find u s.t.

$$\begin{aligned} -\Delta u - k^2 u &= f && \text{in } \Omega, \\ u &= 0 && \text{on } \partial\Omega. \end{aligned}$$

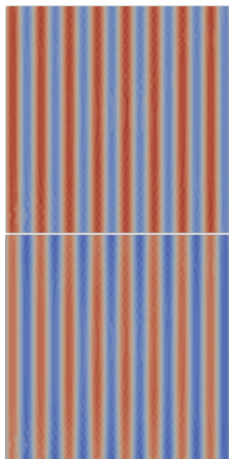
Pollution effect¹: For **fixed** mesh size h , we loose quasi-optimality as the wave-number k **increases**.

Goal: Derive **practical** scheme to generate a mesh that guarantees quasi-optimality

→ implementable

¹extensively studied, e.g. [Babuška, Sauter, 2000], [Melenk, Sauter, 2010 & 2011], [Bernkopf, Chaumont-Frelet, Melenk 2024]

$$u_{\text{ex}} = \sin(16\pi x), \quad k = 5$$



$$\Pi_h^{L^2} u_{\text{ex}} \\ h = 0.075$$

$$u_h \\ h = 0.075$$

Let Ω be **Lipschitz**. For a wave-number k , find u s.t.

$$\begin{aligned} -\Delta u - k^2 u &= f && \text{in } \Omega, \\ u &= 0 && \text{on } \partial\Omega. \end{aligned}$$

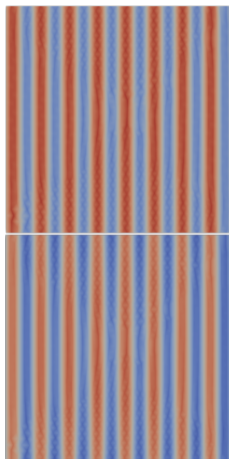
Pollution effect¹: For **fixed** mesh size h , we loose quasi-optimality as the wave-number k **increases**.

Goal: Derive **practical** scheme to generate a mesh that guarantees quasi-optimality

- implementable
- no smoothness assumptions on Ω

¹extensively studied, e.g. [Babuška, Sauter, 2000], [Melenk, Sauter, 2010 & 2011], [Bernkopf, Chaumont-Frelet, Melenk 2024]

$$u_{\text{ex}} = \sin(16\pi x), \quad k = 5$$



$$\Pi_h^{L^2} u_{\text{ex}} \\ h = 0.075$$

$$u_h \\ h = 0.075$$

Theorem (Ciarlet²)

Let X be Hilbert and $a(\cdot, \cdot) : X \times X \rightarrow \mathbb{C}$ be a bounded sesquilinear form. The problem

$$\text{find } u \in X \text{ s.t. } a(u, v) = f(v) \quad \forall v \in X$$

is *well-posed* iff $\exists T : X \rightarrow X$ bijective s.t. $a(\cdot, T\cdot)$ is coercive, i.e.

$$\Re\{a(u, Tu)\} \geq \alpha \|u\|_X^2 \quad \forall u, v \in X.$$

²see e.g., P. Ciarlet Jr., "T-coercivity: Application to the discretization of Helmholtz-like problems", 2012.

Theorem (Ciarlet²)

Let X be Hilbert and $a(\cdot, \cdot) : X \times X \rightarrow \mathbb{C}$ be a bounded sesquilinear form. The problem

$$\text{find } u \in X \text{ s.t. } a(u, v) = f(v) \quad \forall v \in X$$

is *well-posed* iff $\exists T : X \rightarrow X$ bijective s.t. $a(\cdot, T\cdot)$ is coercive, i.e.

$$\Re\{a(u, Tu)\} \geq \alpha \|u\|_X^2 \quad \forall u, v \in X.$$

Note: For Hilbert spaces, T-coercivity is equivalent to the inf sup-condition.

- necessary & sufficient condition for well-posedness
- has to be shown on the discrete level (with uniform constant) to conclude *quasi-optimality*.

²see e.g., P. Ciarlet Jr., "T-coercivity: Application to the discretization of Helmholtz-like problems", 2012.

Let $(\lambda^{(i)}, e^{(i)})_{i \in \mathbb{N}}$ be the eigenpairs of $-\Delta$ on Ω (normed s.t. $\|e^{(i)}\|_{H^1} = 1$). Define

$$i_* = \max\{i \in \mathbb{N} : \lambda^{(i)} < k^2\} \quad (\text{assuming } k^2 \notin \{\lambda^{(i)}\})$$

²P. Ciarlet Jr., "T-coercivity: Application to the discretization of Helmholtz-like problems", 2012.

Let $(\lambda^{(i)}, e^{(i)})_{i \in \mathbb{N}}$ be the eigenpairs of $-\Delta$ on Ω (normed s.t. $\|e^{(i)}\|_{H^1} = 1$). Define

$$i_* = \max\{i \in \mathbb{N} : \lambda^{(i)} < k^2\} \quad (\text{assuming } k^2 \notin \{\lambda^{(i)}\})$$

Then for $u \in H_0^1(\Omega)$

$$\begin{aligned} a(u, u) &:= \int_{\Omega} \nabla u \cdot \nabla u \, dx - k^2 \int_{\Omega} u^2 \, dx \\ &= \sum_{i \leq i_*} \underbrace{\left(\frac{\lambda^{(i)} - k^2}{1 + \lambda^{(i)}} \right)}_{\leq 0} (u^{(i)})^2 + \sum_{i > i_*} \underbrace{\left(\frac{\lambda^{(i)} - k^2}{1 + \lambda^{(i)}} \right)}_{\geq 0} (u^{(i)})^2 \\ &\not\geq \alpha \|u\|_{H^1}^2. \end{aligned}$$

²P. Ciarlet Jr., "T-coercivity: Application to the discretization of Helmholtz-like problems", 2012.

Let $(\lambda^{(i)}, e^{(i)})_{i \in \mathbb{N}}$ be the eigenpairs of $-\Delta$ on Ω (normed s.t. $\|e^{(i)}\|_{H^1} = 1$). Define

$$i_* = \max\{i \in \mathbb{N} : \lambda^{(i)} < k^2\} \quad (\text{assuming } k^2 \notin \{\lambda^{(i)}\})$$

Then for $u \in H_0^1(\Omega)$

$$\begin{aligned} a(u, u) &:= \int_{\Omega} \nabla u \cdot \nabla u \, dx - k^2 \int_{\Omega} u^2 \, dx \\ &= \sum_{i \leq i_*} \underbrace{\left(\frac{\lambda^{(i)} - k^2}{1 + \lambda^{(i)}} \right)}_{\leq 0} (u^{(i)})^2 + \sum_{i > i_*} \underbrace{\left(\frac{\lambda^{(i)} - k^2}{1 + \lambda^{(i)}} \right)}_{\geq 0} (u^{(i)})^2 \\ &\not\geq \alpha \|u\|_{H^1}^2. \end{aligned}$$

→ Construct $T : X \rightarrow X : e^{(i)} \mapsto \begin{cases} -e^{(i)} & \text{if } i \leq i_*, \\ +e^{(i)} & \text{if } i > i_*. \end{cases}$

²P. Ciarlet Jr., "T-coercivity: Application to the discretization of Helmholtz-like problems", 2012.

Let $(\lambda^{(i)}, e^{(i)})_{i \in \mathbb{N}}$ be the eigenpairs of $-\Delta$ on Ω (normed s.t. $\|e^{(i)}\|_{H^1} = 1$). Define

$$i_* = \max\{i \in \mathbb{N} : \lambda^{(i)} < k^2\} \quad (\text{assuming } k^2 \notin \{\lambda^{(i)}\})$$

Then for $u \in H_0^1(\Omega)$

$$\begin{aligned} a(u, Tu) &:= \int_{\Omega} \nabla u \cdot \nabla Tu \, dx - k^2 \int_{\Omega} u Tu \, dx \\ &= \sum_{i \leq i_*} \underbrace{\left(\frac{k^2 - \lambda^{(i)}}{1 + \lambda^{(i)}} \right)}_{\geq 0} (u^{(i)})^2 + \sum_{i > i_*} \underbrace{\left(\frac{\lambda^{(i)} - k^2}{1 + \lambda^{(i)}} \right)}_{\geq 0} (u^{(i)})^2 \\ &\geq \alpha \|u\|_{H^1}^2. \end{aligned}$$

→ Construct $T : X \rightarrow X : e^{(i)} \mapsto \begin{cases} -e^{(i)} & \text{if } i \leq i_*, \\ +e^{(i)} & \text{if } i > i_*. \end{cases} \implies a(\cdot, \cdot)$ is T-coercive.

³P. Ciarlet Jr., "T-coercivity: Application to the discretization of Helmholtz-like problems", 2012.

Let \mathcal{T}_h be a triangulation of Ω with mesh size h and $X_h := \mathbb{P}^p(\mathcal{T}_h) \cap H_0^1(\Omega) \subset H_0^1(\Omega)$, then the Galerkin approximation of the Helmholtz problem reads

$$\text{find } u_h \in X_h \text{ s.t. } a(u_h, v_h) = f(v_h) \quad \forall v_h \in X_h.$$

Let \mathcal{T}_h be a triangulation of Ω with mesh size h and $X_h := \mathbb{P}^p(\mathcal{T}_h) \cap H_0^1(\Omega) \subset H_0^1(\Omega)$, then the Galerkin approximation of the Helmholtz problem reads

$$\text{find } u_h \in X_h \text{ s.t. } a(u_h, v_h) = f(v_h) \quad \forall v_h \in X_h.$$

If $a(\cdot, \cdot)$ is **uniformly** \mathcal{T}_h -coercive, i.e. if $\exists T_h : X_h \rightarrow X_h$ bijective & $\alpha_* > 0$ independent of h s.t.

$$\Re\{a(u_h, T_h u_h)\} \geq \alpha_* \|u_h\|_{H^1}^2 \quad \forall u_h \in X_h,$$

then the H^1 -conforming FEM is **quasi-optimal**.

Let \mathcal{T}_h be a triangulation of Ω with mesh size h and $X_h := \mathbb{P}^p(\mathcal{T}_h) \cap H_0^1(\Omega) \subset H_0^1(\Omega)$, then the Galerkin approximation of the Helmholtz problem reads

$$\text{find } u_h \in X_h \text{ s.t. } a(u_h, v_h) = f(v_h) \quad \forall v_h \in X_h.$$

If $a(\cdot, \cdot)$ is **uniformly** \mathcal{T}_h -coercive, i.e. if $\exists T_h : X_h \rightarrow X_h$ bijective & $\alpha_* > 0$ independent of h s.t.

$$\Re\{a(u_h, T_h u_h)\} \geq \alpha_* \|u_h\|_{H^1}^2 \quad \forall u_h \in X_h,$$

then the H^1 -conforming FEM is **quasi-optimal**.

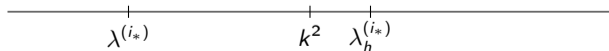
Let $(\lambda_h^{(i)}, e_h^{(i)})_{i \in \mathbb{N}}$ be **conforming** approximations of $(\lambda^{(i)}, e^{(i)})_{i \in \mathbb{N}}$. Then

- $(e_h^{(i)})$ is a basis of X_h
- $\lambda^{(i)} \leq \lambda_h^{(i)}$ for all i , $\lambda_h^{(1)} \leq \lambda_h^{(2)} \leq \dots$

Define $T_h : X_h \rightarrow X_h : e_h^{(i)} \mapsto \begin{cases} -e_h^{(i)} & \text{if } i \leq i_*, \\ +e_h^{(i)} & \text{if } i > i_*. \end{cases}$

For $u_h \in X_h$, we have

$$\begin{aligned} a(u_h, T_h u_h) &= \sum_{i \leq i_*} \underbrace{\left(\frac{k^2 - \lambda_h^{(i)}}{1 + \lambda_h^{(i)}} \right)}_{??} (u_h^{(i)})^2 + \sum_{i > i_*} \underbrace{\left(\frac{\lambda_h^{(i)} - k^2}{1 + \lambda_h^{(i)}} \right)}_{\geq 0} (u_h^{(i)})^2 \\ &\stackrel{??}{\geq} \alpha \|u_h\|_{H^1}^2 \end{aligned}$$



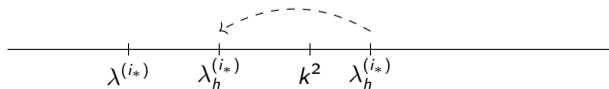
Define $T_h : X_h \rightarrow X_h : e_h^{(i)} \mapsto \begin{cases} -e_h^{(i)} & \text{if } i \leq i_*, \\ +e_h^{(i)} & \text{if } i > i_*. \end{cases}$

For $u_h \in X_h$, we have

$$\begin{aligned} a(u_h, T_h u_h) &= \sum_{i \leq i_*} \underbrace{\left(\frac{k^2 - \lambda_h^{(i)}}{1 + \lambda_h^{(i)}} \right)}_{\geq 0} (u_h^{(i)})^2 + \sum_{i > i_*} \underbrace{\left(\frac{\lambda_h^{(i)} - k^2}{1 + \lambda_h^{(i)}} \right)}_{\geq 0} (u_h^{(i)})^2 \\ &\geq \alpha \|u_h\|_{H^1}^2 \end{aligned}$$

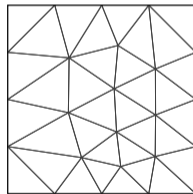
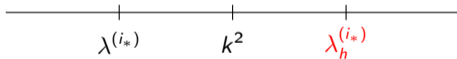
$\rightarrow a(\cdot, \cdot)$ is **uniformly** T_h -coercive **iff** $\lambda_h^{(i_*)} < k^2$.

h small enough

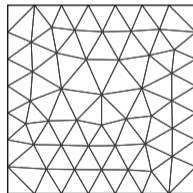
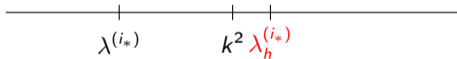


- We can use this criterion to generate a mesh that guarantees quasi-optimality!
1. Determine i_* s.t. $i_* = \max\{i \in \mathbb{N} : \lambda^{(i)} < k^2\}$ (either analytically or numerically).

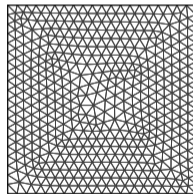
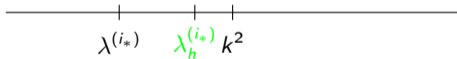
- We can use this criterion to generate a mesh that guarantees quasi-optimality!
1. **Determine** i_* s.t. $i_* = \max\{i \in \mathbb{N} : \lambda^{(i)} < k^2\}$ (either analytically or numerically).
 2. Pick h_0 and solve the Laplace eigenvalue problem. If $\lambda_h^{(i_*)} \geq k^2$, **refine** the mesh and repeat.



- We can use this criterion to generate a mesh that guarantees quasi-optimality!
1. **Determine** i_* s.t. $i_* = \max\{i \in \mathbb{N} : \lambda^{(i)} < k^2\}$ (either analytically or numerically).
 2. Pick h_0 and solve the Laplace eigenvalue problem. If $\lambda_h^{(i_*)} \geq k^2$, **refine** the mesh and repeat.

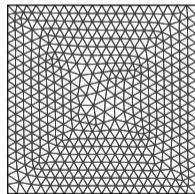
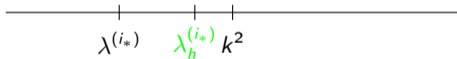


- We can use this criterion to generate a mesh that guarantees quasi-optimality!
1. **Determine** i_* s.t. $i_* = \max\{i \in \mathbb{N} : \lambda^{(i)} < k^2\}$ (either analytically or numerically).
 2. Pick h_0 and solve the Laplace eigenvalue problem. If $\lambda_h^{(i_*)} \geq k^2$, **refine** the mesh and repeat.



→ We can use this criterion to generate a mesh that guarantees quasi-optimality!

1. **Determine** i_* s.t. $i_* = \max\{i \in \mathbb{N} : \lambda^{(i)} < k^2\}$ (either analytically or numerically).
2. Pick h_0 and solve the Laplace eigenvalue problem. If $\lambda_h^{(i_*)} \geq k^2$, **refine** the mesh and repeat.



3. Solve the Helmholtz problem on the mesh obtained in Step 2. Since $\lambda_h^{(i_*)} < k^2$, we have **quasi-optimality**.

Can we optimize the mesh generation process?

→ minimize the number of required mesh elements / dofs

Can we optimize the mesh generation process?

→ minimize the number of required mesh elements / dofs

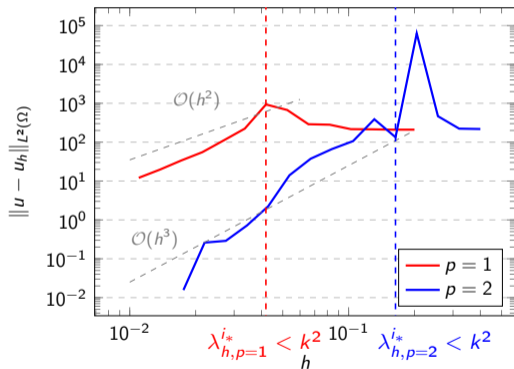
Babuška-Rheinboldt error estimator (averaged over $i_* + \ell$ eigenpairs):

$$\eta = i_*^{-1} \sum_{i=1}^{i_*+\ell} \sum_{K \in \mathcal{T}_h} \left(h_K^2 \|\Delta e_h^{(i)} + \lambda_h^{(i)} e_h^{(i)}\|_{L^2(K)}^2 + \frac{h_K}{2} \|\nabla e_h^{(i)} \cdot n\|_{L^2(\partial K \setminus \partial \Omega)}^2 \right).$$

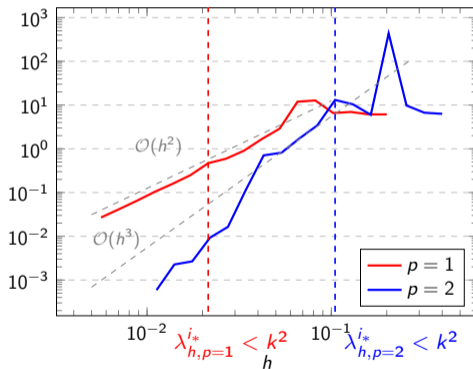
→ use **adaptive refinements** based on this error estimator in Step 2

$\Omega = [0, 1]^2$, Eigenvalues of $-\Delta$: $\lambda_{i,j} = \pi^2(i^2 + j^2)$, $i, j \in \mathbb{N}$.

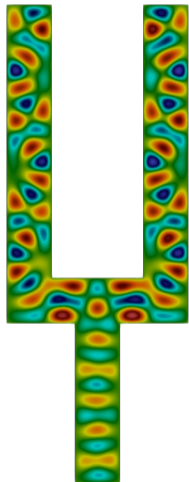
$k = 10$



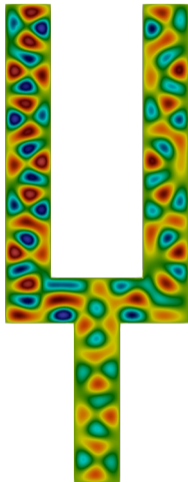
$k = 20$



Numerical Example - Tuning fork ($k = 10$)

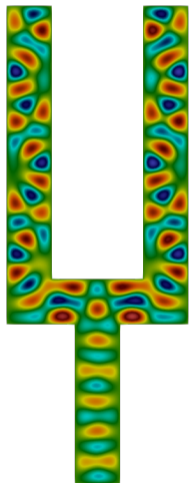


$$\lambda_h^{(i_*)} > k^2$$

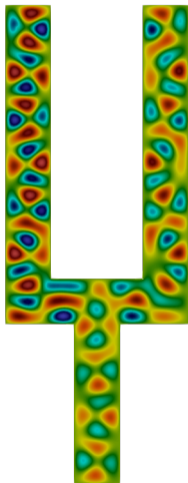


$$\lambda_h^{(i_*)} < k^2$$

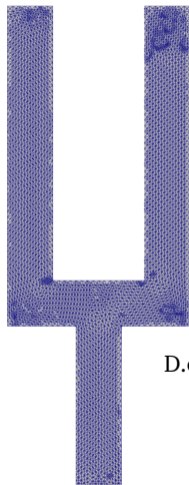
Numerical Example - Tuning fork ($k = 10$)



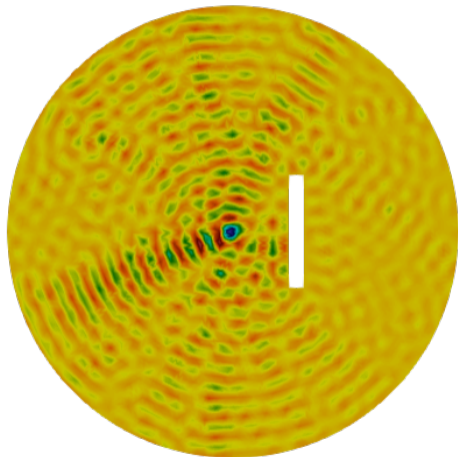
$$\lambda_h^{(i_*)} > k^2$$



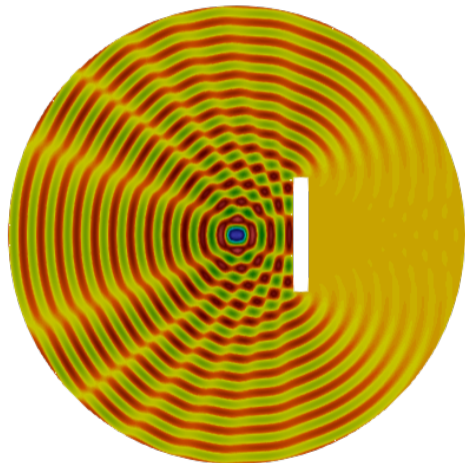
$$\lambda_h^{(i_*)} < k^2$$



D.o.f.s. required s.t. $\lambda_h^{(i_*)} < k^2$:
uniform: **336,449**
adaptive: **161,102**



$$\lambda_h^{(i_*)} > k^2$$



$$\lambda_h^{(i_*)} < k^2$$


Conclusions

- quasi-optimality is intimately connected to the discrete eigenvalues of $-\Delta$
- we can use this connection to generate a mesh that guarantees quasi-optimality:
 - determine maximal index i_* s.t. $\lambda^{(i_*)} < k^2$
 - adaptively refine the mesh until $\lambda_h^{(i_*)} < k^2$
 - Solve the Helmholtz problem
- can be extended to Robin / Mixed boundary conditions

Conclusions

- quasi-optimality is intimately connected to the discrete eigenvalues of $-\Delta$
- we can use this connection to generate a mesh that guarantees quasi-optimality:
 - determine maximal index i_* s.t. $\lambda^{(i_*)} < k^2$
 - adaptively refine the mesh until $\lambda_h^{(i_*)} < k^2$
 - Solve the Helmholtz problem
- can be extended to Robin / Mixed boundary conditions

Curious to learn more?

 TvB, U. Zerbinati, "An adaptive mesh refinement strategy to ensure quasi-optimality of the conforming finite element method for the Helmholtz equation via T-coercivity" (2024), <https://arxiv.org/pdf/2403.06266>.

