

Modelling Round-off Error in the Fast Gradient Method for Predictive Control

Ian McInerney¹, Eric C. Kerrigan^{1,2}, George A. Constantinides¹

¹Department of Electrical and Electronic Engineering, Imperial College London

²Department of Aeronautics, Imperial College London

Optimal Control Problem

- Constrained Linear Quadratic Regulator

$$\min_{u,x} \frac{1}{2} x'_N P x_N + \frac{1}{2} \sum_{k=0}^{N-1} \begin{bmatrix} x_k \\ u_k \end{bmatrix}' \begin{bmatrix} Q & S \\ S' & R \end{bmatrix} \begin{bmatrix} x_k \\ u_k \end{bmatrix}$$

$$\text{s.t. } x_{k+1} = Ax_k + Bu_k, \quad k = 0, \dots, N-1$$

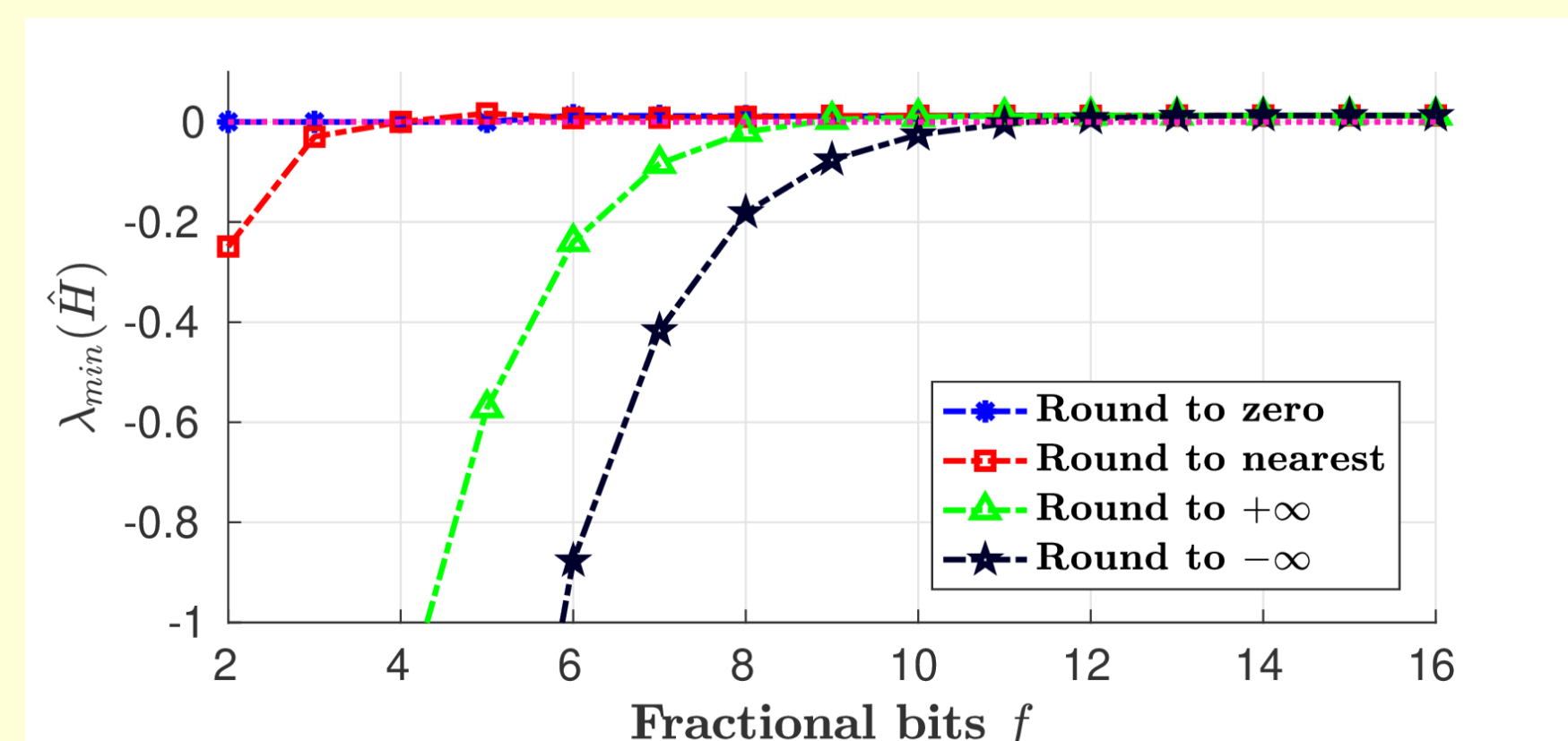
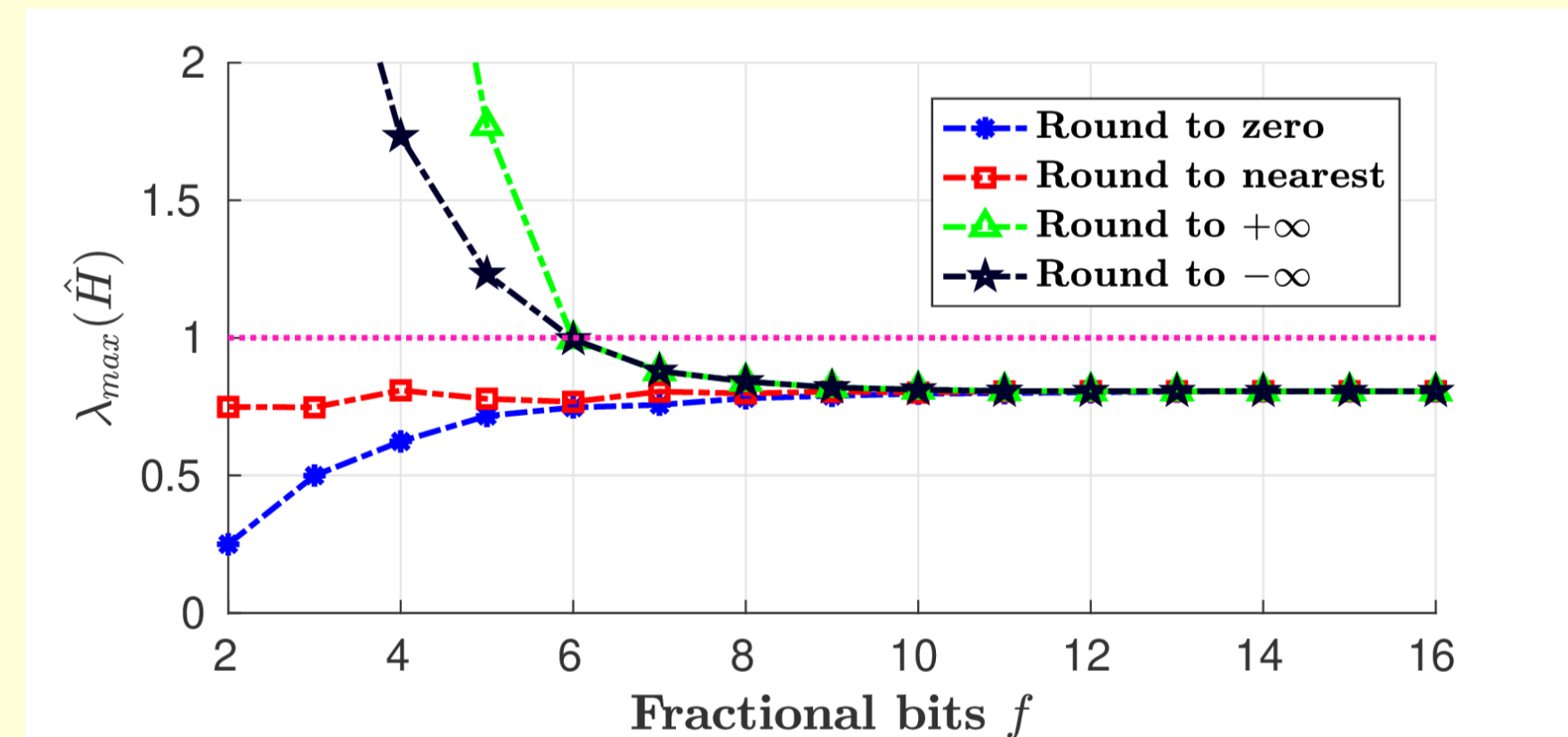
$$Fu_k \leq c_u, \quad k = 0, \dots, N-1$$

Remove state variables

$$\min_u \frac{1}{2} u' H_c u + x'_0 J' u$$

$$\text{s.t. } Gu \leq F \hat{x}_0 + g$$

- Solve using the Fast Gradient Method
 - For stability and convexity: $\lambda(H_c) \in (0, 1)$
- Converting problem to fixed-point representation can cause loss of convexity and stability



Rounding Model

- Model rounding loss as a perturbation matrix, and define the largest allowable perturbation

Definition 1 (Rounding stability margin).

Let $\hat{H} = H + E$ with $\|E\|_2 = \beta$ and $\lambda(H) \in (0, 1)$. The rounding stability margin η is the smallest value of β that causes the eigenvalues of \hat{H} to leave the interval $(0, 1)$.

- Use the pseudospectrum of H_c to compute η
 - For generic matrices, add $\pm \epsilon$ to all terms as the perturbation, giving

$$\eta(H) = \min \left\{ \|(-H)^{-1}\|_2^{-1}, \|(I-H)^{-1}\|_2^{-1} \right\}.$$
 - For Schur-stable systems, exploit the Toeplitz structure of H_c , and use its matrix symbol, giving

$$\eta(H) = \min \left\{ \|(-\mathcal{P}_H)^{-1}\|_{H_\infty}^{-1}, \|(I_m - \mathcal{P}_H)^{-1}\|_{H_\infty}^{-1} \right\},$$

$$\mathcal{P}_{H_c P}(z) := (z\mathcal{G}(z)_s)^* Q (z\mathcal{G}(z)_s) + R \quad \forall z \in \{z \in \mathbb{C} : |z| = 1\}.$$

- For a system with m inputs and a horizon of length N , size the fractional bits as follows:

– For generic matrices:

$$f = \begin{cases} \lceil -\log_2 \left(\frac{\eta}{mN} \right) \rceil - 1 & \text{if using round to nearest,} \\ \lceil -\log_2 \left(\frac{\eta}{mN} \right) \rceil & \text{otherwise.} \end{cases}$$

– For Schur-stable predicted systems using round to nearest or round to zero, solve:

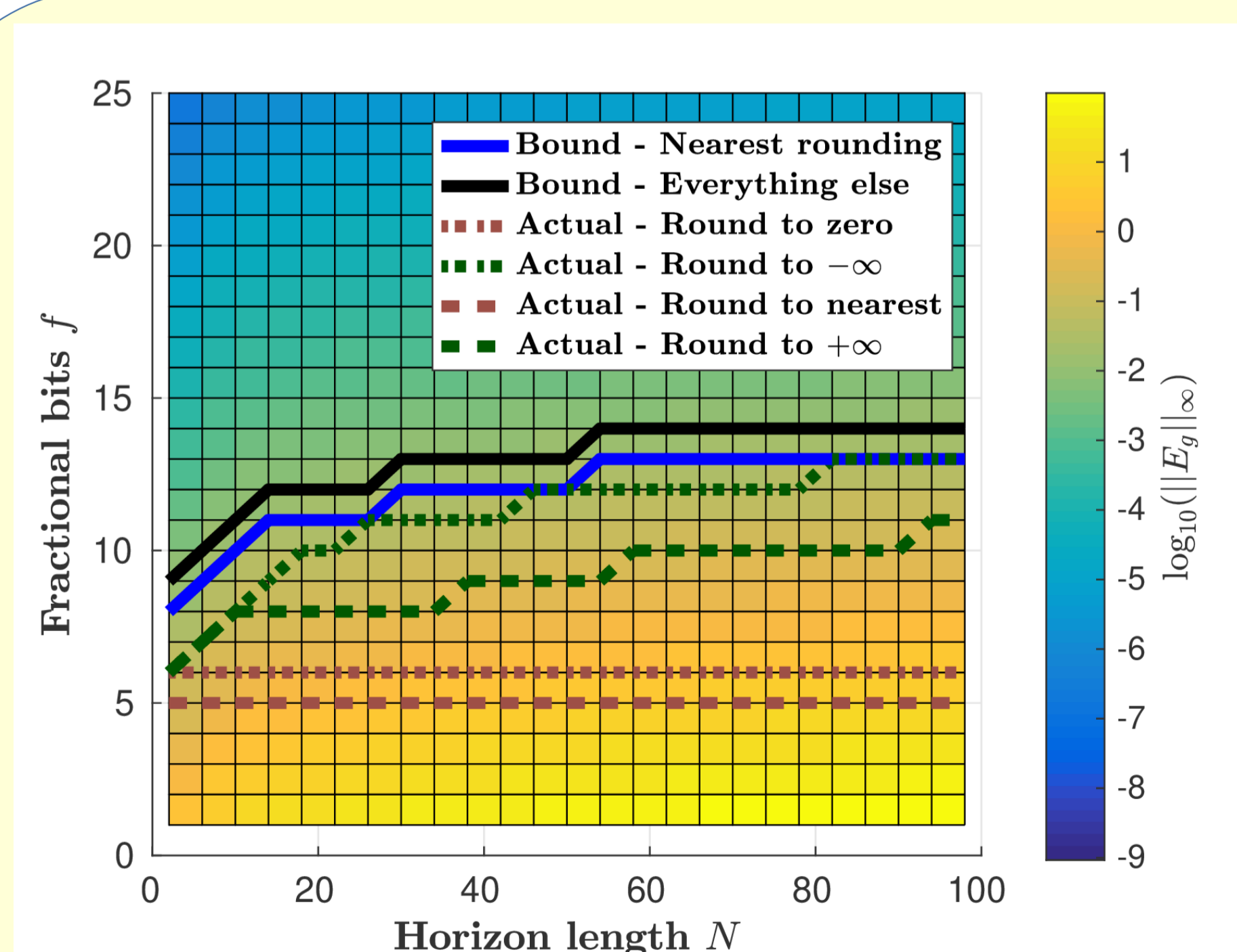
$$\min f$$

$$\text{s.t. } |\epsilon_f| m(2k-1) + 2 \|\mathcal{P}_{\hat{H}}(k, \cdot)\|_{H_\infty} < \eta$$

where

$$\mathcal{G}_P := \begin{cases} x^+ = Ax + Bu \\ y = B'Px \end{cases}, \quad \mathcal{P}_n(z) := \sum_{i=0}^{n-1} A^i z^{-i} \quad \forall z \in \{z \in \mathbb{C} : |z| = 1\}$$

$$\mathcal{P}_{\hat{H}}(n, z) := z\mathcal{G}_P(z) - B'PP_n(z)B \quad \forall z \in \{z \in \mathbb{C} : |z| = 1\}.$$



Minimum fractional bits needed for a given horizon length

- Exploiting the Toeplitz structure leads to a reduction of the number of fractional bits by 40% compared to generic matrices.
- The minimum fractional length uses 77% less memory, 33% fewer DSP blocks, and is 25% faster compared with the floating-point version

Fractional Length	Logic Resources				Power (mW)	Solve Time (μ s)
	LUT	FF	DSP	BRAM		
f=12	947	768	4	2	20	532.17
f=16	1,136	912	4	2	25	612.17
f=21	887	1,033	8	8	43	701.77
f=26	993	1,237	12	9	48	701.77
Float (single)	2,161	1,545	5	14	51	982.17

Resource usage for the Fast Gradient Method on a Zynq 7020