

Absolute error in complex fixed-point arithmetic

September 11, 2015

This draft proves a few results on the loss of precision that occurs when doing basic complex operations. Throughout this part, we note $z_j = x_j + iy_j$ a complex number and $\tilde{z}_j = \tilde{x}_j + i\tilde{y}_j$ its approximation.

1 General case

We suppose that

$$|x_j - \tilde{x}_j| \leq k_{R,j}2^{-P} \quad |y_j - \tilde{y}_j| \leq k_{I,j}2^{-P}$$

I don't know if this section is really likely to happen (different error on the real and imaginary parts, really?) But it's what they do in the MPC document...

1.0.1 Addition

If $z = z_1 + z_2$ we have

$$\begin{aligned} |x_1 + x_2 - (\tilde{x}_1 + \tilde{x}_2)| &\leq |x_1 - \tilde{x}_1| + |x_2 - \tilde{x}_2| \leq (k_{R,1} + k_{R,2})2^{-P} \\ |y_1 + y_2 - (\tilde{y}_1 + \tilde{y}_2)| &\leq |y_1 - \tilde{y}_1| + |y_2 - \tilde{y}_2| \leq (k_{I,1} + k_{I,2})2^{-P} \end{aligned}$$

So "the errors add up".

1.0.2 Multiplication

If $z = z_1 z_2$ we have $x = x_1 x_2 - y_1 y_2$ and $y = x_1 y_2 + x_2 y_1$. We write

$$\begin{aligned} |x_1 x_2 - \tilde{x}_1 \tilde{x}_2| &\leq \frac{1}{2} |(x_1 - \tilde{x}_1)(x_2 + \tilde{x}_2) + (x_1 + \tilde{x}_1)(x_2 - \tilde{x}_2)| \\ &\leq \frac{1}{2} (|x_1 - \tilde{x}_1|(|x_2| + |\tilde{x}_2|) + (|x_1| + |\tilde{x}_1|)|x_2 - \tilde{x}_2|) \\ &\leq \frac{1}{2} (k_{R,1}(|x_2| + |\tilde{x}_2|) + k_{R,2}(|x_1| + |\tilde{x}_1|))2^{-P} \\ &\leq (k_{R,1}|x_2| + k_{R,2}|x_1|)2^{-P} + k_{R,1}2^{-P}k_{R,2}2^{-P} \end{aligned}$$

We suppose that $k_{R,1} \leq 2^{P/2}$ and $k_{R,2} \leq 2^{P/2}$, which means that the majority of the bits are correct (if not, we're in trouble); hence

$$|x_1 x_2 - \tilde{x}_1 \tilde{x}_2| \leq (k_{R,1}|x_2| + k_{R,2}|x_1| + 1)2^{-P}$$

Similarly we have

$$\begin{aligned} |y_1 y_2 - \tilde{y}_1 \tilde{y}_2| &\leq (k_{I,1}|y_2| + k_{I,2}|y_1| + 1)2^{-P} \\ |x_1 y_2 - \tilde{x}_1 \tilde{y}_2| &\leq (k_{R,1}|y_2| + k_{I,2}|x_1| + 1)2^{-P} \\ |y_1 x_2 - \tilde{y}_1 \tilde{x}_2| &\leq (k_{I,1}|x_2| + k_{R,2}|y_1| + 1)2^{-P} \end{aligned}$$

Hence

$$\begin{aligned} |x - \tilde{x}| &\leq (2 + k_{R,1}|x_2| + k_{R,2}|x_1| + k_{I,1}|y_2| + k_{I,2}|y_1|)2^{-P} \\ |y - \tilde{y}| &\leq (2 + k_{R,1}|y_2| + k_{I,2}|x_1| + k_{I,1}|x_2| + k_{R,2}|y_1|)2^{-P} \end{aligned}$$

Bounding the norms of real and imaginary parts by the norm of the number itself:

$$\begin{aligned} |x - \tilde{x}| &\leq (2 + (k_{R,1} + k_{I,1})|z_2| + (k_{R,2} + k_{I,2})|z_1|)2^{-P} \\ |y - \tilde{y}| &\leq (2 + (k_{R,1} + k_{I,1})|z_2| + (k_{R,2} + k_{I,2})|z_1|)2^{-P} \end{aligned}$$

1.0.3 Squaring

The above formulas can be simplified in the case $z = z_1^2$:

$$\begin{aligned} |x_1^2 - \tilde{x}_1^2| &\leq (2k_{R,1}|x_1| + 1)2^{-P} \\ |y_1^2 - \tilde{y}_1^2| &\leq (2k_{I,1}|y_1| + 1)2^{-P} \\ |x_1 y_1 - \tilde{x}_1 \tilde{y}_1| &\leq (k_{R,1}|y_1| + k_{I,1}|x_1| + 1)2^{-P} \end{aligned}$$

which gives the bounds

$$\begin{aligned} |x - \tilde{x}| &\leq (2 + 2k_{R,1}|x_1| + 2k_{I,1}|y_1|)2^{-P} \\ |y - \tilde{y}| &\leq (2 + 2k_{R,1}|y_1| + 2k_{I,1}|x_1|)2^{-P} \end{aligned}$$

Bounding the norms of real and imaginary parts by the norm of the number itself gives:

$$\begin{aligned} |x - \tilde{x}| &\leq (2 + 2(k_{R,1} + k_{I,1})|z_1|)2^{-P} \\ |y - \tilde{y}| &\leq (2 + 2(k_{R,1} + k_{I,1})|z_1|)2^{-P} \end{aligned}$$

1.0.4 Norm

We wish to compute $z = |z_1|^2$. We use previous estimates to write:

$$\begin{aligned} ||z|^2 - |\tilde{z}|^2| &\leq |x_1^2 - \tilde{x}_1^2| + |y_1^2 - \tilde{y}_1^2| \\ &\leq (4 + 2k_{R,1}|x_1| + 2k_{I,1}|y_1| + 2k_{R,1}|y_1| + 2k_{I,1}|x_1|)2^{-P} \end{aligned}$$

Again, bounding the norms of real and imaginary parts by the norm of the number itself:

$$||z|^2 - |\tilde{z}|^2| \leq (4 + 4(k_{R,1} + k_{I,1})|z_1|)2^{-P}$$

1.0.5 Division of a real by a positive real

Let us assume $|a_1 - \tilde{a}_1| \leq k_1 2^{-P}$ and $|a_2 - \tilde{a}_2| \leq k_2 2^{-P}$, with $a_2 > 0$. In addition, assume that $a_2 > k_2 2^{-P}$; this means that the sign of a_2 is known and that we are not at risk of dividing by 0. Then we have

$$\begin{aligned}
\left| \frac{a_1}{a_2} - \frac{\tilde{a}_1}{\tilde{a}_2} \right| &\leq \left| \frac{a_1 \tilde{a}_2 - a_2 \tilde{a}_1}{a_2 \tilde{a}_2} \right| \\
&\leq \left| \frac{(a_1 - \tilde{a}_1)(a_2 + \tilde{a}_2) - (a_1 a_2 - \tilde{a}_1 \tilde{a}_2)}{a_2 \tilde{a}_2} \right| \\
&\leq \frac{a_2 + \tilde{a}_2}{a_2 \tilde{a}_2} k_1 2^{-P} + \frac{k_1 |a_2| + k_2 |a_1| + 1}{a_2 \tilde{a}_2} 2^{-P} \\
&\leq \left(\frac{2a_2 + \tilde{a}_2}{a_2 \tilde{a}_2} k_1 + \frac{k_2 |a_1| + 1}{a_2 \tilde{a}_2} \right) 2^{-P} \\
&\leq \left(\frac{3a_2 + k_2 2^{-P}}{a_2 \tilde{a}_2} k_1 + \frac{k_2 |a_1| + 1}{a_2 \tilde{a}_2} \right) 2^{-P} \\
&\leq \left(\frac{3k_1}{a_2 - k_2 2^{-P}} + \frac{k_2 (|a_1| + k_1 2^{-P}) + 1}{a_2 (a_2 - k_2 2^{-P})} \right) 2^{-P}
\end{aligned}$$

To further simplify, we assume that $a_2 \geq 2k_2 2^{-P}$; otherwise, we might end up in a case where not even the high bit of \tilde{a}_2 is correct (for instance if $a_2 = 2k_2 2^{-P}$ and $\tilde{a}_2 = k_2 2^{-P}$). This means that $a_2 - k_2 2^{-P} \geq a_2/2$ and helps with denominators:

$$\left| \frac{a_1}{a_2} - \frac{\tilde{a}_1}{\tilde{a}_2} \right| \leq \left(\frac{6k_1}{a_2} + \frac{2k_2 (|a_1| + k_1 2^{-P}) + 2}{a_2^2} \right) 2^{-P}$$

Making the same supposition for a_1 further simplifies:

$$\left| \frac{a_1}{a_2} - \frac{\tilde{a}_1}{\tilde{a}_2} \right| \leq \left(\frac{6k_1}{a_2} + \frac{4k_2 |a_1| + 2}{a_2^2} \right) 2^{-P}$$

1.0.6 Complex division

We write $\frac{z_1}{z_2} = \frac{z_1 \bar{z}_2}{|z_2|^2}$ and then chain the results. So:

$$\begin{aligned}
||z_2|^2 - |\tilde{z}_2|^2| &\leq (4 + 2k_{R,2}|x_2| + 2k_{I,2}|y_2| + 2k_{R,2}|y_2| + 2k_{I,2}|x_2|) 2^{-P} \\
|\operatorname{Re}(z_1 \bar{z}_2) - \operatorname{Re}(\tilde{z}_1 \bar{\tilde{z}}_2)| &\leq (2 + k_{R,1}|x_2| + k_{R,2}|x_1| + k_{I,1}|y_2| + k_{I,2}|y_1|) 2^{-P} \\
|\operatorname{Im}(z_1 \bar{z}_2) - \operatorname{Im}(\tilde{z}_1 \bar{\tilde{z}}_2)| &\leq (2 + k_{R,1}|y_2| + k_{I,2}|x_1| + k_{I,1}|x_2| + k_{R,2}|y_1|) 2^{-P}
\end{aligned}$$

and we use the previous part:

$$\begin{aligned}
\left| \operatorname{Re} \left(\frac{z_1}{z_2} \right) - \operatorname{Re} \left(\frac{\tilde{z}_1}{\tilde{z}_2} \right) \right| &\leq \left(\frac{3(2 + k_{R,1}|x_2| + k_{R,2}|x_1| + k_{I,1}|y_2| + k_{I,2}|y_1|)}{|z_2|^2 - (4 + 2k_{R,2}|x_2| + 2k_{I,2}|y_2| + 2k_{R,2}|y_2| + 2k_{I,2}|x_2|)2^{-P}} \right. \\
&\quad \left. + (4 + 2k_{R,2}|x_2| + 2k_{I,2}|y_2| + 2k_{R,2}|y_2| + 2k_{I,2}|x_2|) \times \right. \\
&\quad \left. \frac{(|\operatorname{Re}(\tilde{z}_1 \tilde{z}_2)| + (2 + k_{R,1}|x_2| + k_{R,2}|x_1| + k_{I,1}|y_2| + k_{I,2}|y_1|)2^{-P})}{|z_2|^2(|z_2|^2 - (4 + 2k_{R,2}|x_2| + 2k_{I,2}|y_2| + 2k_{R,2}|y_2| + 2k_{I,2}|x_2|)2^{-P})} + \right. \\
&\quad \left. \frac{1}{|z_2|^2(|z_2|^2 - (4 + 2k_{R,2}|x_2| + 2k_{I,2}|y_2| + 2k_{R,2}|y_2| + 2k_{I,2}|x_2|)2^{-P})} \right) 2^{-P} \\
\left| \operatorname{Im} \left(\frac{z_1}{z_2} \right) - \operatorname{Im} \left(\frac{\tilde{z}_1}{\tilde{z}_2} \right) \right| &\leq \left(\frac{3(2 + k_{R,1}|y_2| + k_{I,2}|x_1| + k_{I,1}|x_2| + k_{R,2}|y_1|)}{|z_2|^2 - (4 + 2k_{R,2}|x_2| + 2k_{I,2}|y_2| + 2k_{R,2}|y_2| + 2k_{I,2}|x_2|)2^{-P}} \right. \\
&\quad \left. + (4 + 2k_{R,2}|x_2| + 2k_{I,2}|y_2| + 2k_{R,2}|y_2| + 2k_{I,2}|x_2|) \times \right. \\
&\quad \left. \frac{(|\operatorname{Im}(\tilde{z}_1 \tilde{z}_2)| + (2 + k_{R,1}|y_2| + k_{I,2}|x_1| + k_{I,1}|x_2| + k_{R,2}|y_1|)2^{-P})}{|z_2|^2(|z_2|^2 - (4 + 2k_{R,2}|x_2| + 2k_{I,2}|y_2| + 2k_{R,2}|y_2| + 2k_{I,2}|x_2|)2^{-P})} \right. \\
&\quad \left. + \frac{1}{|z_2|^2(|z_2|^2 - (4 + 2k_{R,2}|x_2| + 2k_{I,2}|y_2| + 2k_{R,2}|y_2| + 2k_{I,2}|x_2|)2^{-P})} \right) 2^{-P}
\end{aligned}$$

Let's make the supposition that

$$|z_2|^2 \geq 2 \times (4 + 2k_{R,2}|x_2| + 2k_{I,2}|y_2| + 2k_{R,2}|y_2| + 2k_{I,2}|x_2|)2^{-P}$$

to simplify the denominators:

$$\begin{aligned}
\left| \operatorname{Re} \left(\frac{z_1}{z_2} \right) - \operatorname{Re} \left(\frac{\tilde{z}_1}{\tilde{z}_2} \right) \right| &\leq \left(\frac{6(2 + k_{R,1}|x_2| + k_{R,2}|x_1| + k_{I,1}|y_2| + k_{I,2}|y_1|)}{|z_2|^2} \right. \\
&\quad + 2(4 + 2k_{R,2}|x_2| + 2k_{I,2}|y_2| + 2k_{R,2}|y_2| + 2k_{I,2}|x_2|) \times \\
&\quad \left. \frac{(|\operatorname{Re}(\tilde{z}_1 \tilde{z}_2)| + (2 + k_{R,1}|x_2| + k_{R,2}|x_1| + k_{I,1}|y_2| + k_{I,2}|y_1|)2^{-P})}{|z_2|^4} + \frac{2}{|z_2|^4} \right) 2^{-P} \\
&\leq \left(\frac{6(2 + k_{R,1}|x_2| + k_{R,2}|x_1| + k_{I,1}|y_2| + k_{I,2}|y_1|)}{|z_2|^2} \right. \\
&\quad + 2(4 + 2k_{R,2}|x_2| + 2k_{I,2}|y_2| + 2k_{R,2}|y_2| + 2k_{I,2}|x_2|) \times \\
&\quad \left. \frac{(|\tilde{x}_1 \tilde{x}_2| + |\tilde{y}_1 \tilde{y}_2|) + (2 + k_{R,1}|x_2| + k_{R,2}|x_1| + k_{I,1}|y_2| + k_{I,2}|y_1|)2^{-P}}{|z_2|^4} + \frac{2}{|z_2|^4} \right) 2^{-P} \\
&\leq \left(\frac{6(2 + k_{R,1}|x_2| + k_{R,2}|x_1| + k_{I,1}|y_2| + k_{I,2}|y_1|)}{|z_2|^2} \right. \\
&\quad + 2(4 + 2k_{R,2}|x_2| + 2k_{I,2}|y_2| + 2k_{R,2}|y_2| + 2k_{I,2}|x_2|) \times \\
&\quad \left. \frac{|x_1||x_2| + |y_1||y_2| + 2(1 + k_{R,1}|x_2| + k_{R,2}|x_1| + k_{I,1}|y_2| + k_{I,2}|y_1|)2^{-P} + (k_{R,1}k_{R,2} + k_{I,1}k_{I,2})}{|z_2|^4} \right) 2^{-P} \\
\left| \operatorname{Im} \left(\frac{z_1}{z_2} \right) - \operatorname{Im} \left(\frac{\tilde{z}_1}{\tilde{z}_2} \right) \right| &\leq \left(\frac{6(2 + k_{R,1}|y_2| + k_{I,2}|x_1| + k_{I,1}|x_2| + k_{R,2}|y_1|)}{|z_2|^2} \right. \\
&\quad + 2(4 + 2k_{R,2}|x_2| + 2k_{I,2}|y_2| + 2k_{R,2}|y_2| + 2k_{I,2}|x_2|) \times \\
&\quad \left. \frac{(|\operatorname{Im}(\tilde{z}_1 \tilde{z}_2)| + (2 + k_{R,1}|y_2| + k_{I,2}|x_1| + k_{I,1}|x_2| + k_{R,2}|y_1|)2^{-P})}{|z_2|^4} + \frac{2}{|z_2|^4} \right) 2^{-P} \\
&\leq \left(\frac{6(2 + k_{R,1}|y_2| + k_{I,2}|x_1| + k_{I,1}|x_2| + k_{R,2}|y_1|)}{|z_2|^2} \right. \\
&\quad + 2(4 + 2k_{R,2}|x_2| + 2k_{I,2}|y_2| + 2k_{R,2}|y_2| + 2k_{I,2}|x_2|) \times \\
&\quad \left. \frac{(|\tilde{x}_1 \tilde{y}_2| + |\tilde{x}_2 \tilde{y}_1|) + (2 + k_{R,1}|y_2| + k_{I,2}|x_1| + k_{I,1}|x_2| + k_{R,2}|y_1|)2^{-P}}{|z_2|^4} + \frac{2}{|z_2|^4} \right) 2^{-P} \\
&\leq \left(\frac{6(2 + k_{R,1}|y_2| + k_{I,2}|x_1| + k_{I,1}|x_2| + k_{R,2}|y_1|)}{|z_2|^2} \right. \\
&\quad + 2(4 + 2k_{R,2}|x_2| + 2k_{I,2}|y_2| + 2k_{R,2}|y_2| + 2k_{I,2}|x_2|) \times \\
&\quad \left. \frac{(|x_1||y_2| + |x_2||y_1|) + 2(1 + k_{R,1}|y_2| + k_{I,2}|x_1| + k_{I,1}|x_2| + k_{R,2}|y_1|)2^{-P} + (k_{R,1}k_{I,2} + k_{I,1}k_{R,2})}{|z_2|^4} \right) 2^{-P}
\end{aligned}$$

Suppose that

$$\begin{aligned}
2(1 + k_{R,1}|y_2| + k_{I,2}|x_1| + k_{I,1}|x_2| + k_{R,2}|y_1|)2^{-P} + (k_{R,1}k_{I,2} + k_{I,1}k_{R,2})2^{-2P} &\leq 1 \\
2(1 + k_{R,1}|x_2| + k_{R,2}|x_1| + k_{I,1}|y_2| + k_{I,2}|y_1|)2^{-P} + (k_{R,1}k_{R,2} + k_{I,1}k_{I,2})2^{-2P} &\leq 1
\end{aligned}$$

which is true for instance if $k_1 \leq 2^{P/2}$ and $|x_1| \leq 2^{P/2-3}$ (and the same for the others). We can then simplify:

$$\begin{aligned} \left| \operatorname{Re} \left(\frac{z_1}{z_2} \right) - \operatorname{Re} \left(\frac{\tilde{z}_1}{\tilde{z}_2} \right) \right| &\leq \left(\frac{6(2 + k_{R,1}|x_2| + k_{R,2}|x_1| + k_{I,1}|y_2| + k_{I,2}|y_1|)}{|z_2|^2} \right. \\ &\quad \left. + \frac{2(4 + 2k_{R,2}|x_2| + 2k_{I,2}|y_2| + 2k_{R,2}|y_2| + 2k_{I,2}|x_2|)(|x_1||x_2| + |y_1||y_2| + 1) + 2}{|z_2|^4} \right) 2^{-P} \\ \left| \operatorname{Im} \left(\frac{z_1}{z_2} \right) - \operatorname{Im} \left(\frac{\tilde{z}_1}{\tilde{z}_2} \right) \right| &\leq \left(\frac{6(2 + k_{R,1}|y_2| + k_{I,2}|x_1| + k_{I,1}|x_2| + k_{R,2}|y_1|)}{|z_2|^2} \right. \\ &\quad \left. + \frac{2(4 + 2k_{R,2}|x_2| + 2k_{I,2}|y_2| + 2k_{R,2}|y_2| + 2k_{I,2}|x_2|)(|x_1||y_2| + |x_2||y_1| + 1) + 2}{|z_2|^4} \right) 2^{-P} \end{aligned}$$

Finally, we simplify using $|x_1|, |y_1| \leq |z_1|$ and $|x_2|, |y_2| \leq |z_2|$:

$$\begin{aligned} \left| \operatorname{Re} \left(\frac{z_1}{z_2} \right) - \operatorname{Re} \left(\frac{\tilde{z}_1}{\tilde{z}_2} \right) \right| &\leq \left(\frac{6(2 + (k_{R,1} + k_{I,1})|z_2| + (k_{R,2} + k_{I,2})|z_1|)}{|z_2|^2} \right. \\ &\quad \left. + \frac{2(4 + 4(k_{R,2} + k_{I,2})|z_2|)(2|z_1||z_2| + 1) + 2}{|z_2|^4} \right) 2^{-P} \\ \left| \operatorname{Im} \left(\frac{z_1}{z_2} \right) - \operatorname{Im} \left(\frac{\tilde{z}_1}{\tilde{z}_2} \right) \right| &\leq \left(\frac{6(2 + (k_{R,1} + k_{I,1})|z_2| + (k_{R,2} + k_{I,2})|z_1|)}{|z_2|^2} \right. \\ &\quad \left. + \frac{2(4 + 4(k_{R,2} + k_{I,2})|z_2|)(2|z_1||z_2| + 1) + 2}{|z_2|^4} \right) 2^{-P} \end{aligned}$$

2 Case where $k_R = k_I$

In this section, we assume that

$$|x_j - \tilde{x}_j| \leq k_j 2^{-P} \quad |y_j - \tilde{y}_j| \leq k_j 2^{-P}$$

The results simplify as follows:

2.1 Addition

$$\begin{aligned} |x_1 + x_2 - (\tilde{x}_1 + \tilde{x}_2)| &\leq (k_1 + k_2) 2^{-P} \\ |y_1 + y_2 - (\tilde{y}_1 + \tilde{y}_2)| &\leq (k_1 + k_2) 2^{-P} \end{aligned}$$

2.2 Multiplication

$$\begin{aligned} |x - \tilde{x}| &\leq (2 + 2k_1|z_2| + 2k_2|z_1|) 2^{-P} \\ |y - \tilde{y}| &\leq (2 + 2k_1|z_2| + 2k_2|z_1|) 2^{-P} \end{aligned}$$

2.3 Squaring

$$\begin{aligned} |x - \tilde{x}| &\leq (2 + 4k_1|z_1|) 2^{-P} \\ |y - \tilde{y}| &\leq (2 + 4k_1|z_1|) 2^{-P} \end{aligned}$$

2.4 Norm

$$||z|^2 - |\tilde{z}|^2| \leq (8k_1|z_1| + 4)2^{-P}$$

2.5 Complex division

$$\begin{aligned} \left| \operatorname{Re} \left(\frac{z_1}{z_2} \right) - \operatorname{Re} \left(\frac{\tilde{z}_1}{\tilde{z}_2} \right) \right| &\leq \left(\frac{6(2 + 2k_1|z_2| + 2k_2|z_1|)}{|z_2|^2} + \frac{2(4 + 8k_2|z_2|)(2|z_1||z_2| + 1) + 2}{|z_2|^4} \right) 2^{-P} \\ \left| \operatorname{Im} \left(\frac{z_1}{z_2} \right) - \operatorname{Im} \left(\frac{\tilde{z}_1}{\tilde{z}_2} \right) \right| &\leq \left(\frac{6(2 + 2k_1|z_2| + 2k_2|z_1|)}{|z_2|^2} + \frac{2(4 + 8k_2|z_2|)(2|z_1||z_2| + 1) + 2}{|z_2|^4} \right) 2^{-P} \end{aligned}$$

2.6 Square root

We write

$$\begin{aligned} |\sqrt{z_1} - \sqrt{\tilde{z}_1}| &= \frac{|z_1 - \tilde{z}_1|}{|\sqrt{z_1} + \sqrt{\tilde{z}_1}|} \\ &\leq \frac{\sqrt{2}k_1}{|\sqrt{z_1} + \sqrt{\tilde{z}_1}|} 2^{-P} \end{aligned}$$

If we suppose that z_1 and \tilde{z}_1 are in the same quadrant, which is true if we suppose $|x_1| > k_1 2^{-P}$ and $|y_1| > k_1 2^{-P}$, then $\sqrt{z_1}$ and $\sqrt{\tilde{z}_1}$ are in the same quadrant (since the angle is just divided by 2). This means that $|\sqrt{z_1} + \sqrt{\tilde{z}_1}| \geq |\sqrt{z_1}|$. Hence

$$|\sqrt{z_1} - \sqrt{\tilde{z}_1}| \leq \frac{\sqrt{2}k_1}{\sqrt{|z_1|}} 2^{-P}$$

2.7 Exponential

Starting with real numbers: we have $|e^x - e^{\tilde{x}}| \leq e^t |x - \tilde{x}|$ with $t \in]x, \tilde{x}[$ by Taylor-Lagrange with order 1 / Rolle's theorem. Hence

$$\begin{aligned} |e^x - e^{\tilde{x}}| &\leq e^t |x - \tilde{x}| \\ &\leq \max(e^x, e^{\tilde{x}}) k_x 2^{-P} \\ &\leq e^x e^{k_x 2^{-P}} k_x 2^{-P} \end{aligned}$$

Since $k_x 2^{-P} < 1/2$ (and we usually suppose $k_x 2^{-P} \leq 2^{-P/2}$ so that at least half of the bits are correct), we have $e^x \leq 1 + x + \frac{x^2}{2} \frac{1}{1-x} \leq 1 + x + x^2 \leq 1 + 2x$. Hence

$$\begin{aligned} |e^x - e^{\tilde{x}}| &\leq e^x (1 + 2k_x 2^{-P}) k_x 2^{-P} \\ &\leq e^x (k_x + 2) 2^{-P} \end{aligned}$$

Now for complex numbers:

$$\begin{aligned} |e^{x+iy} - e^{\tilde{x}+i\tilde{y}}| &\leq |e^x - e^{\tilde{x}+i(\tilde{y}-y)}| \\ &\leq \sqrt{(e^x - e^{\tilde{x}} \cos(\tilde{y} - y))^2 + \sin^2(\tilde{y} - y)} \end{aligned}$$

Since for positive numbers $a + b \leq a + b + 2\sqrt{a}\sqrt{b}$, $\sqrt{a+b} \leq \sqrt{a} + \sqrt{b}$ and

$$\begin{aligned} |e^{x+iy} - e^{\tilde{x}+i\tilde{y}}| &\leq |e^x - e^{\tilde{x}} \cos(\tilde{y} - y)| + e^{\tilde{y}} |\sin(\tilde{y} - y)| \\ &\leq |e^x - e^{\tilde{x}}| + e^{\tilde{x}} (|1 - \cos(\tilde{y} - y)| + |\sin(\tilde{y} - y)|) \\ &\leq e^x (k_x + 2) 2^{-P} + e^x (1 + 2k_x 2^{-P}) (|1 - \cos(\tilde{y} - y)| + |\sin(\tilde{y} - y)|) \end{aligned}$$

For $x > 0$ we have $\sin(x) \leq x$ and $|1 - \cos(x)| \leq \frac{x^2}{2}$ (since $\cos x = 1 - 2\sin^2(x/2)$ or the theorem for alternate series), hence

$$\begin{aligned} |e^{x+iy} - e^{\tilde{x}+i\tilde{y}}| &\leq e^x (k_x + 2) 2^{-P} + e^x (1 + 2k_x 2^{-P}) ((k_y 2^{-P})^2 / 4 + k_y 2^{-P}) \\ &\leq (e^x (k_x + 2)) 2^{-P} + e^x (1 + 2k_x 2^{-P}) (1/4 + k_y 2^{-P}) \\ &\leq e^x (5/2 k_x + k_y + 4.25) 2^{-P} \end{aligned}$$

because as always we suppose $k_x 2^{-P} \leq 2^{-P/2}$. If we suppose $k_x = k_y$, then the bound is

$$|e^z - e^{\tilde{z}}| \leq e^x \frac{7k_x + 8.5}{2} 2^{-P}$$

3 Comparison with the MPC document

Remember that we have:

$$\epsilon^+ = \frac{\max(0, x - \tilde{x})}{|\tilde{x}|} \leq \frac{|x - \tilde{x}|}{\tilde{x}}, \quad \epsilon^- = \frac{\max(0, \tilde{x} - x)}{|\tilde{x}|} \leq \frac{|x - \tilde{x}|}{\tilde{x}}$$

Also we have $Exp(x) = \lceil \log_2 |x| \rceil + 1$, so that $2^{Exp(x)-1} \leq |x|$. The error bounds are:

$$\text{error}(x_1) \leq k_1 2^{Exp(x_1)-p}$$

We take their final results and see if they're the same as ours. Spoiler: probably due to approximations they do (most notably bounding $|x|$ by $2^{Exp(x)}$), their bounds are worse; so we can't really say our work directly derives from the MPC document, but our proofs are the same spirit and techniques for sure.

3.1 Addition

They write: $\text{error}(x) \leq k_{R,1} 2^{Exp(\tilde{x}_1)-p} + k_{R,2} 2^{Exp(\tilde{x}_2)-p}$, and so

$$\begin{aligned} \text{error}(x) &\leq k_{R,1} 2^{Exp(\tilde{x}_1)-p} + k_{R,2} 2^{Exp(\tilde{x}_2)-p} \\ &\leq \text{error}(\tilde{x}_1) + \text{error}(\tilde{x}_2) \end{aligned}$$

Our bound was $\text{error}(\tilde{x}_1) + \text{error}(\tilde{x}_2)$, but with $\text{error}(\tilde{x}_1) \leq k_{R,1} 2^{-P}$.

3.2 Multiplication

Let's see:

$$\begin{aligned} |x_1 x_2 - \tilde{x}_1 \tilde{x}_2| &\leq (k_{R,1} (2 + \epsilon_{R,2}^+) + k_{R,2} (2 + \epsilon_{R,1}^+)) 2^{Exp(\tilde{x}_1 \tilde{x}_2)-p} \\ &\leq k_{R,1} (2 + \epsilon_{R,2}^+) 2^{Exp(\tilde{x}_1) + Exp(\tilde{x}_2) - p} + k_{R,2} (2 + \epsilon_{R,1}^+) 2^{Exp(\tilde{x}_1) + Exp(\tilde{x}_2) - p} \\ &\leq \text{error}(\tilde{x}_1) (2^{Exp(\tilde{x}_2)+1} + \epsilon_{R,2}^+ 2^{Exp(\tilde{x}_2)}) + \text{error}(\tilde{x}_2) (2^{Exp(\tilde{x}_1)+1} + \epsilon_{R,1}^+ 2^{Exp(\tilde{x}_1)}) \\ &\leq \text{error}(\tilde{x}_1) (2|x_2| + 2|x_2 - \tilde{x}_2|) + \text{error}(\tilde{x}_2) (2|x_1| + 2|x_1 - \tilde{x}_1|) \end{aligned}$$

Now suppose that $\text{error}(\tilde{x}_1)|x_2 - \tilde{x}_2|$ is smaller than 2^{-P} ; this is kinda like what we do when saying $k_1 k_2 \leq 2^P$. And note: the MPC document does the same thing, i.e. between Equations 8 and 9: "Under normal circumstances $\epsilon_1 \epsilon_2 \leq 2^{1-P}$ ". Then we have

$$|x_1 x_2 - \tilde{x}_1 \tilde{x}_2| \leq 2(|x_1| \text{error}(\tilde{x}_2) + |x_2| \text{error}(\tilde{x}_1) + 1)$$

and this is worse by a factor 2 than our bound.

Note: That supposition is not valid in the case given in the article ! Since in the article we say that "at most we will lose $17P$ so we should work at precision $18P$ ". It's not true : if you want that to be true, you should work at precision $34P$! Otherwise the "1" in the formulas will be even bigger than that, e.g. 2^{16P} . The naive algorithm is probably immune to that because $\log B$ is fine, but our quasi-optimal algorithm is really not; we should work at really big precision then, and that's two times asymptotically worse.

3.3 Norm

They write:

$$\begin{aligned} |x_1 - \tilde{x}_1|^2 &\leq 2(k_{R,1}(2 + \epsilon_{R,1}^+) + k_{I,1}(2 + \epsilon_{I,1}^+))2^{Exp(\tilde{x}-p)} \\ &\leq 8|x_1^2|k_{R,1}2^{-P} + 4k_{R,1}|x_1 - \tilde{x}_1||\tilde{x}_1|2^{-P} \end{aligned}$$

Interestingly, this is worse for $|x_1| > 1$, but better if not (because $|z|^2$ instead of $|z|$).

Their other expression

$$\text{error}(\tilde{x}) = 2k_1(2 + \epsilon_1)2^{Exp(\tilde{x})-p}$$

is the same deal.

3.4 Division

They don't even do till the end, handwaving some details

3.5 Sqrt

Interesting:

$$\left| \frac{\sqrt{x} - \sqrt{\tilde{x}}}{\sqrt{\tilde{x}}} \right| \leq \frac{\epsilon_1}{2\sqrt{1 - \epsilon_1}}$$

with $\epsilon_1 = \left| \frac{x - \tilde{x}}{\tilde{x}} \right|$. So

$$\begin{aligned} \left| \frac{\sqrt{x} - \sqrt{\tilde{x}}}{\sqrt{\tilde{x}}} \right| &\leq \frac{\epsilon_1}{2\sqrt{1 - \epsilon_1}} \\ &\leq \frac{|x - \tilde{x}|}{2|\sqrt{x}|\sqrt{|x| - |x - \tilde{x}|}} \\ |\sqrt{x} - \sqrt{\tilde{x}}| &\leq \frac{|x - \tilde{x}|}{2\sqrt{|x| - |x - \tilde{x}|}} \end{aligned}$$

Suppose that $|x| - |x - \tilde{x}| \geq |x|/2$ and you find our bound. So we didn't do too bad.

4 What should probably go in the final pdf / appendix of the paper

We outline here bounds on the absolute error that can arise when computing elementary operations in fixed point arithmetic with numbers in precision P . We write $z_k = x_k + iy_k$, and $\tilde{z}_k = \tilde{x}_k + i\tilde{y}_k$ its approximation with fixed precision P . We adopt the following convention:

$$|x_j - \tilde{x}_j| \leq k_j 2^{-P} \quad |y_j - \tilde{y}_j| \leq k_j 2^{-P}$$

that is to say, the approximation of a complex number z_j has real and imaginary parts bounded by k_j times the smallest gap; hence we have $|z_j - \tilde{z}_j| \leq \sqrt{2}k_j 2^{-P}$.

4.1 Addition

If $z = z_1 + z_2$ we have

$$\begin{aligned} |x_1 + x_2 - (\tilde{x}_1 + \tilde{x}_2)| &\leq |x_1 - \tilde{x}_1| + |x_2 - \tilde{x}_2| \leq (k_1 + k_2)2^{-P} \\ |y_1 + y_2 - (\tilde{y}_1 + \tilde{y}_2)| &\leq |y_1 - \tilde{y}_1| + |y_2 - \tilde{y}_2| \leq (k_1 + k_2)2^{-P} \end{aligned}$$

So "the errors add up".

4.2 Multiplication

If $z = z_1 z_2$ we have $x = x_1 x_2 - y_1 y_2$ and $y = x_1 y_2 + x_2 y_1$. We write

$$\begin{aligned} |x_1 x_2 - \tilde{x}_1 \tilde{x}_2| &\leq \frac{1}{2} |(x_1 - \tilde{x}_1)(x_2 + \tilde{x}_2) + (x_1 + \tilde{x}_1)(x_2 - \tilde{x}_2)| \\ &\leq \frac{1}{2} (|x_1 - \tilde{x}_1|(|x_2| + |\tilde{x}_2|) + (|x_1| + |\tilde{x}_1|)|x_2 - \tilde{x}_2|) \\ &\leq \frac{1}{2} (k_1(|x_2| + |\tilde{x}_2|) + k_2(|x_1| + |\tilde{x}_1|)) 2^{-P} \\ &\leq (k_1|x_2| + k_2|x_1|)2^{-P} + k_1 2^{-P} k_2 2^{-P} \end{aligned}$$

We suppose that $k_1 \leq 2^{P/2}$ and $k_2 \leq 2^{P/2}$, which means that the majority of the bits are correct (if not, we're in trouble); hence

$$|x_1 x_2 - \tilde{x}_1 \tilde{x}_2| \leq (k_1|x_2| + k_2|x_1| + 1)2^{-P}$$

Similarly we have

$$\begin{aligned} |y_1 y_2 - \tilde{y}_1 \tilde{y}_2| &\leq (k_1|y_2| + k_2|y_1| + 1)2^{-P} \\ |x_1 y_2 - \tilde{x}_1 \tilde{y}_2| &\leq (k_1|y_2| + k_2|x_1| + 1)2^{-P} \\ |y_1 x_2 - \tilde{y}_1 \tilde{x}_2| &\leq (k_1|x_2| + k_2|y_1| + 1)2^{-P} \end{aligned}$$

Hence

$$\begin{aligned} |x - \tilde{x}| &\leq (2 + k_1|x_2| + k_2|x_1| + k_1|y_2| + k_2|y_1|)2^{-P} \\ |y - \tilde{y}| &\leq (2 + k_1|y_2| + k_2|x_1| + k_1|x_2| + k_2|y_1|)2^{-P} \end{aligned}$$

Bounding the norms of real and imaginary parts by the norm of the number itself:

$$\begin{aligned} |x - \tilde{x}| &\leq (2 + 2k_1|z_2| + 2k_2|z_1|)2^{-P} \\ |y - \tilde{y}| &\leq (2 + 2k_1|z_2| + 2k_2|z_1|)2^{-P} \end{aligned}$$

4.3 Squaring

The above formulas can be simplified in the case $z = z_1^2$:

$$\begin{aligned} |x_1^2 - \tilde{x}_1^2| &\leq (2k_1|x_1| + 1)2^{-P} \\ |y_1^2 - \tilde{y}_1^2| &\leq (2k_1|y_1| + 1)2^{-P} \\ |x_1y_1 - \tilde{x}_1\tilde{y}_1| &\leq (k_1|y_1| + k_1|x_1| + 1)2^{-P} \end{aligned}$$

which gives the bounds

$$\begin{aligned} |x - \tilde{x}| &\leq (2 + 2k_1(|x_1| + |y_1|))2^{-P} \\ |y - \tilde{y}| &\leq (2 + 2k_1(|x_1| + |y_1|))2^{-P} \end{aligned}$$

Bounding the norms of real and imaginary parts by the norm of the number itself gives:

$$\begin{aligned} |x - \tilde{x}| &\leq (2 + 4k_1|z_1|)2^{-P} \\ |y - \tilde{y}| &\leq (2 + 4k_1|z_1|)2^{-P} \end{aligned}$$

4.4 Norm

We wish to compute $z = |z_1|^2$. We use previous estimates to write:

$$\begin{aligned} ||z|^2 - |\tilde{z}|^2| &\leq |x_1^2 - \tilde{x}_1^2| + |y_1^2 - \tilde{y}_1^2| \\ &\leq (4 + 4k_1|x_1| + 4k_1|y_1|)2^{-P} \end{aligned}$$

Again, bounding the norms of real and imaginary parts by the norm of the number itself:

$$||z|^2 - |\tilde{z}|^2| \leq (4 + 8k_1|z_1|)2^{-P}$$

4.5 Division of a real by a positive real

Let us assume $|a_1 - \tilde{a}_1| \leq k_12^{-P}$ and $|a_2 - \tilde{a}_2| \leq k_22^{-P}$, with $a_2 > 0$. In addition, assume that $a_2 > k_22^{-P}$; this means that the sign of a_2 is known and that we are not at risk of dividing by 0. Then we have

$$\begin{aligned} \left| \frac{a_1}{a_2} - \frac{\tilde{a}_1}{\tilde{a}_2} \right| &\leq \left| \frac{a_1\tilde{a}_2 - a_2\tilde{a}_1}{a_2\tilde{a}_2} \right| \\ &\leq \left| \frac{(a_1 - \tilde{a}_1)(a_2 + \tilde{a}_2) - (a_1a_2 - \tilde{a}_1\tilde{a}_2)}{a_2\tilde{a}_2} \right| \\ &\leq \frac{a_2 + \tilde{a}_2}{a_2\tilde{a}_2}k_12^{-P} + \frac{k_1|a_2| + k_2|a_1| + 1}{a_2\tilde{a}_2}2^{-P} \\ &\leq \left(\frac{2a_2 + \tilde{a}_2}{a_2\tilde{a}_2}k_1 + \frac{k_2|a_1| + 1}{a_2\tilde{a}_2} \right) 2^{-P} \\ &\leq \left(\frac{3a_2 + k_22^{-P}}{a_2\tilde{a}_2}k_1 + \frac{k_2|a_1| + 1}{a_2\tilde{a}_2} \right) 2^{-P} \\ &\leq \left(\frac{3k_1}{a_2 - k_22^{-P}} + \frac{k_2(|a_1| + k_12^{-P}) + 1}{a_2(a_2 - k_22^{-P})} \right) 2^{-P} \end{aligned}$$

To further simplify, we assume that $a_2 \geq 2k_22^{-P}$; otherwise, we might end up in a case where not even the high bit of \tilde{a}_2 is correct (for instance if $a_2 = 2k_22^{-P}$ and $\tilde{a}_2 = k_22^{-P}$). This means that $a_2 - k_22^{-P} \geq a/2$ and helps with denominators:

$$\left| \frac{a_1}{a_2} - \frac{\tilde{a}_1}{\tilde{a}_2} \right| \leq \left(\frac{6k_1}{a_2} + \frac{2k_2(|a_1| + k_12^{-P}) + 2}{a_2^2} \right) 2^{-P}$$

Making the same supposition for a_1 further simplifies:

$$\left| \frac{a_1}{a_2} - \frac{\tilde{a}_1}{\tilde{a}_2} \right| \leq \left(\frac{6k_1}{a_2} + \frac{4k_2|a_1| + 2}{a_2^2} \right) 2^{-P}$$

4.6 Complex division

We write $\frac{z_1}{z_2} = \frac{z_1\bar{z}_2}{|z_2|^2}$ and then chain the results. To simplify, we bound $|x_i|, |y_i|$ by the norm of the number itself right away. Hence:

$$\begin{aligned} ||z_2|^2 - |\tilde{z}_2|^2| &\leq (4 + 8k_2|z_2|)2^{-P} \\ |\operatorname{Re}(z_1\bar{z}_2) - \operatorname{Re}(\tilde{z}_1\tilde{z}_2)| &\leq (2 + 2k_1|z_2| + 2k_2|z_1|)2^{-P} \\ |\operatorname{Im}(z_1\bar{z}_2) - \operatorname{Im}(\tilde{z}_1\tilde{z}_2)| &\leq (2 + 2k_1|z_2| + 2k_2|z_1|)2^{-P} \end{aligned}$$

and we use the previous part:

$$\begin{aligned}
\left| \operatorname{Re} \left(\frac{z_1}{z_2} \right) - \operatorname{Re} \left(\frac{\tilde{z}_1}{\tilde{z}_2} \right) \right| &\leq \left(\frac{6(2 + 2k_1|z_2| + 2k_2|z_1|)}{|z_2|^2} \right. \\
&\quad \left. + 2(4 + 8k_2|z_2|) \times \frac{(|\operatorname{Re}(\tilde{z}_1 \tilde{z}_2)| + (2 + 2k_1|z_2| + 2k_2|z_1|)2^{-P})}{|z_2|^4} + \frac{2}{|z_2|^4} \right) 2^{-P} \\
&\leq \left(\frac{6(2 + 2k_1|z_2| + 2k_2|z_1|)}{|z_2|^2} \right. \\
&\quad \left. + 2(4 + 8k_2|z_2|) \times \frac{(|\tilde{x}_1 \tilde{x}_2| + |\tilde{y}_1 \tilde{y}_2|) + (2 + 2k_1|z_2| + 2k_2|z_1|)2^{-P}}{|z_2|^4} + \frac{2}{|z_2|^4} \right) 2^{-P} \\
&\leq \left(\frac{6(2 + 2k_1|z_2| + 2k_2|z_1|)}{|z_2|^2} \right. \\
&\quad \left. + 2(4 + 8k_2|z_2|) \times \frac{2|z_1||z_2| + 2(1 + 2k_1|z_2| + 2k_2|z_1|)2^{-P} + 2k_1k_22^{-2P}}{|z_2|^4} + \frac{2}{|z_2|^4} \right) 2^{-P} \\
&\quad \text{(using } |\tilde{x}_1| \leq |x_1| + k_12^{-P} \leq |z_1| + k_12^{-P} \text{)}
\end{aligned}$$

$$\begin{aligned}
\left| \operatorname{Im} \left(\frac{z_1}{z_2} \right) - \operatorname{Im} \left(\frac{\tilde{z}_1}{\tilde{z}_2} \right) \right| &\leq \left(\frac{6(2 + 2k_1|z_2| + 2k_2|z_1|)}{|z_2|^2} \right. \\
&\quad \left. + 2(4 + 8k_2|z_2|) \times \frac{(|\operatorname{Im}(\tilde{z}_1 \tilde{z}_2)| + (2 + 2k_1|z_2| + 2k_2|z_1|)2^{-P})}{|z_2|^4} + \frac{2}{|z_2|^4} \right) 2^{-P} \\
&\leq \left(\frac{6(2 + 2k_1|z_2| + 2k_2|z_1|)}{|z_2|^2} \right. \\
&\quad \left. + 2(4 + 8k_2|z_2|) \times \frac{(|\tilde{x}_1 \tilde{y}_2| + |\tilde{x}_2 \tilde{y}_1|) + (2 + 2k_1|z_2| + 2k_2|z_1|)2^{-P}}{|z_2|^4} + \frac{2}{|z_2|^4} \right) 2^{-P} \\
&\leq \left(\frac{6(2 + 2k_1|z_2| + 2k_2|z_1|)}{|z_2|^2} \right. \\
&\quad \left. + 2(4 + 8k_2|z_2|) \times \frac{2|z_1||z_2| + 2(1 + 2k_1|z_2| + 2k_2|z_1|)2^{-P} + 2k_1k_22^{-2P}}{|z_2|^4} + \frac{2}{|z_2|^4} \right) 2^{-P}
\end{aligned}$$

Suppose that

$$\begin{aligned}
2(1 + 2k_1|z_2| + 2k_2|z_1|)2^{-P} + 2k_1k_22^{-2P} &\leq 1 \\
2(1 + 2k_1|z_2| + 2k_2|z_1|)2^{-P} + 2k_1k_22^{-2P} &\leq 1
\end{aligned}$$

which is true for instance if $k_1 \leq 2^{P/2}$ and $|z_1| \leq 2^{P/2-3}$ (and the same for the others). We can then simplify:

$$\begin{aligned}
\left| \operatorname{Re} \left(\frac{z_1}{z_2} \right) - \operatorname{Re} \left(\frac{\tilde{z}_1}{\tilde{z}_2} \right) \right| &\leq \left(\frac{6(2 + 2k_1|z_2| + 2k_2|z_1|)}{|z_2|^2} + \frac{2(4 + 8k_2|z_2|)(2|z_1||z_2| + 1) + 2}{|z_2|^4} \right) 2^{-P} \\
\left| \operatorname{Im} \left(\frac{z_1}{z_2} \right) - \operatorname{Im} \left(\frac{\tilde{z}_1}{\tilde{z}_2} \right) \right| &\leq \left(\frac{6(2 + 2k_1|z_2| + 2k_2|z_1|)}{|z_2|^2} + \frac{2(4 + 8k_2|z_2|)(2|z_1||z_2| + 1) + 2}{|z_2|^4} \right) 2^{-P}
\end{aligned}$$

4.7 Square root

We write

$$\begin{aligned} |\sqrt{z_1} - \sqrt{\tilde{z}_1}| &= \frac{|z_1 - \tilde{z}_1|}{|\sqrt{z_1} + \sqrt{\tilde{z}_1}|} \\ &\leq \frac{\sqrt{2}k_1}{|\sqrt{z_1} + \sqrt{\tilde{z}_1}|} 2^{-P} \end{aligned}$$

If we suppose that z_1 and \tilde{z}_1 are in the same quadrant, which is true if we suppose $|x_1| > k_1 2^{-P}$ and $|y_1| > k_1 2^{-P}$, then $\sqrt{z_1}$ and $\sqrt{\tilde{z}_1}$ are in the same quadrant (since the angle is just divided by 2). This means that $|\sqrt{z_1} + \sqrt{\tilde{z}_1}| \geq |\sqrt{z_1}|$. Hence

$$|\sqrt{z_1} - \sqrt{\tilde{z}_1}| \leq \frac{\sqrt{2}k_1}{\sqrt{|z_1|}} 2^{-P}$$

4.8 Exponential

Starting with real numbers: we have $|e^x - e^{\tilde{x}}| \leq e^t |x - \tilde{x}|$ with $t \in]x, \tilde{x}[$ by Taylor-Lagrange with order 1 / Rolle's theorem. Hence

$$\begin{aligned} |e^x - e^{\tilde{x}}| &\leq e^t |x - \tilde{x}| \\ &\leq \max(e^x, e^{\tilde{x}}) k_x 2^{-P} \\ &\leq e^x e^{k_x 2^{-P}} k_x 2^{-P} \end{aligned}$$

Since $k_x 2^{-P} < 1/2$ (and we usually suppose $k_x 2^{-P} \leq 2^{-P/2}$ so that at least half of the bits are correct), we have $e^x \leq 1 + x + \frac{x^2}{2} \frac{1}{1-x} \leq 1 + x + x^2 \leq 1 + 2x$. Hence

$$\begin{aligned} |e^x - e^{\tilde{x}}| &\leq e^x (1 + 2k_x 2^{-P}) k_x 2^{-P} \\ &\leq e^x (k_x + 2) 2^{-P} \end{aligned}$$

Now for complex numbers:

$$\begin{aligned} |e^{x+iy} - e^{\tilde{x}+i\tilde{y}}| &\leq |e^x - e^{\tilde{x}+i(\tilde{y}-y)}| \\ &\leq \sqrt{(e^x - e^{\tilde{x}} \cos(\tilde{y} - y))^2 + \sin^2(\tilde{y} - y)} \end{aligned}$$

Since for positive numbers $a + b \leq a + b + 2\sqrt{a}\sqrt{b}$, $\sqrt{a+b} \leq \sqrt{a} + \sqrt{b}$ and

$$\begin{aligned} |e^{x+iy} - e^{\tilde{x}+i\tilde{y}}| &\leq |e^x - e^{\tilde{x}} \cos(\tilde{y} - y)| + e^{\tilde{y}} |\sin(\tilde{y} - y)| \\ &\leq |e^x - e^{\tilde{x}}| + e^{\tilde{x}} (|1 - \cos(\tilde{y} - y)| + |\sin(\tilde{y} - y)|) \\ &\leq e^x (k_x + 2) 2^{-P} + e^x (1 + 2k_x 2^{-P}) (|1 - \cos(\tilde{y} - y)| + |\sin(\tilde{y} - y)|) \end{aligned}$$

For $x > 0$ we have $\sin(x) \leq x$ and $|1 - \cos(x)| \leq \frac{x^2}{2}$ (since $\cos x = 1 - 2\sin^2(x/2)$ or the theorem for alternate series), hence

$$\begin{aligned} |e^{x+iy} - e^{\tilde{x}+i\tilde{y}}| &\leq e^x (k_x + 2) 2^{-P} + e^x (1 + 2k_x 2^{-P}) ((k_y 2^{-P})^2 / 4 + k_y 2^{-P}) \\ &\leq (e^x (k_x + 2)) 2^{-P} + e^x (1 + 2k_x 2^{-P}) (1/4 + k_y 2^{-P}) \\ &\leq e^x (5/2 k_x + k_y + 4.25) 2^{-P} \end{aligned}$$

because as always we suppose $k_x 2^{-P} \leq 2^{-P/2}$. If we suppose $k_x = k_y$, then the bound is

$$|e^z - e^{\tilde{z}}| \leq e^x \frac{7k_x + 8.5}{2} 2^{-P}$$