**1.** *Given trees $T_1$ and $T_2$ and matching threshold $r$, to ensure a pq-Gram distance $d$ such that $0 \le d \le r \le 1$, then we must have $|I_1| \ge \frac{1-r}{1+r}|I_2|$.*

*Proof.* Without loss of generality, assume $|T_1| \le |T_2|$. Let $I_1$ be the pq-Gram profile corresponding to $T_1$, and let $I_2$ be the pq-Gram profile corresponding to $T_2$.

Recall that the pq-Gram distance is defined as

$$dist(T_1, T_2) = 1 - 2\frac{|I_1 \cap I_2|}{|I_1 \uplus I_2|}.$$

Note that in any case $|I_1 \uplus I_2| = |I_1| + |I_2|$, and in the best case scenario, to maximize $|I_1 \cap I_2|$, $I_1 \subseteq I_2$. Then $|I_1 \cap I_2| = |I_1|$, so we can make the simplification

$$dist(T_1, T_2) = 1 - 2\frac{|I_1 \cap I_2|}{|I_1 \uplus I_2|} = 1 - 2\frac{|I_1|}{|I_1| + |I_2|}.$$

Then if $r$ is the threshold such that $0 \le r \le 1$, we want distance $d$ to be such that $0 \le d \le r$, so we must have

$$
\begin{aligned}
r &\ge 1 - 2\frac{|I_1|}{|I_1| + |I_2|} \\
\Rightarrow \quad 2\frac{|I_1|}{|I_1| + |I_2|} &\ge 1 - r \\
\Rightarrow \quad \frac{2}{1-r}|I_1| &\ge |I_1| + |I_2| \\
\Rightarrow \quad \frac{2}{1-r}|I_1| - |I_1| &\ge |I_2| \\
\Rightarrow \quad \frac{2-(1-r)}{1-r}|I_1| &\ge |I_2| \\
\Rightarrow \quad \frac{1+r}{1-r}|I_1| &\ge |I_2| \\
\Rightarrow \quad |I_1| &\ge \frac{1-r}{1+r}|I_2|.
\end{aligned}
$$

$\square$