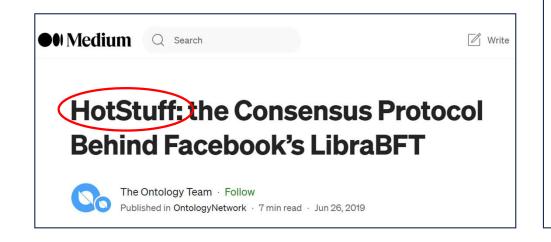
Y.C. Tay

(keynote for IWSF&SHIFT 2025, Sao Paolo, Brazil)

Y.C. Tay

(keynote for IWSF&SHIFT 2025, São Paulo, Brazil)

objective: analytical model that expresses performance in terms of protocol and topology parameters



The Istanbul BFT Consensus Algorithm

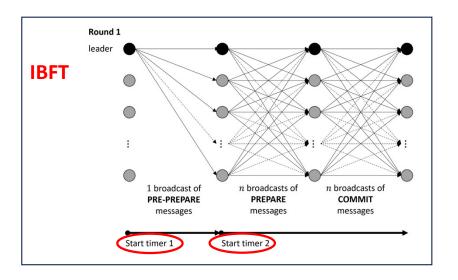
Henrique Moniz

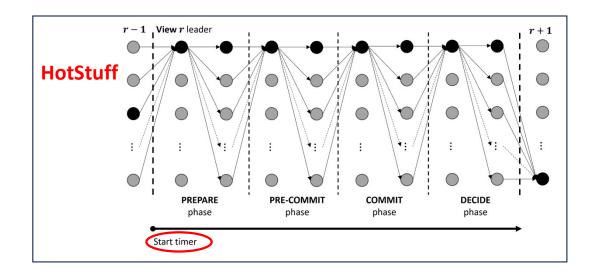
Quorum Engineering

May 20, 2020

Abstract

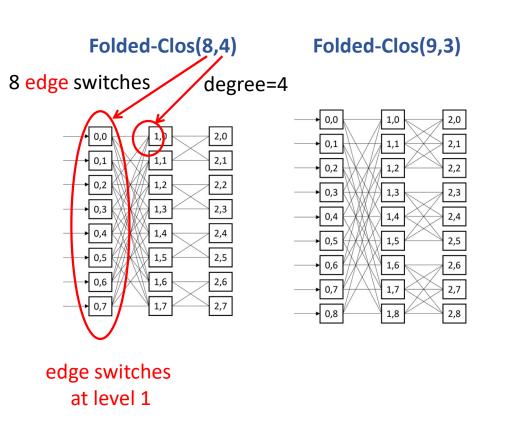
This paper presents IBFT, a simple and elegant Byzantine fault-tolerant consensus algorithm that is used to implement state machine replication in the *Quorum* blockchain. IBFT assumes a partially synchronous communication model, where safety does not depend on any timing assumptions and only liveness depends on periods of synchrony.

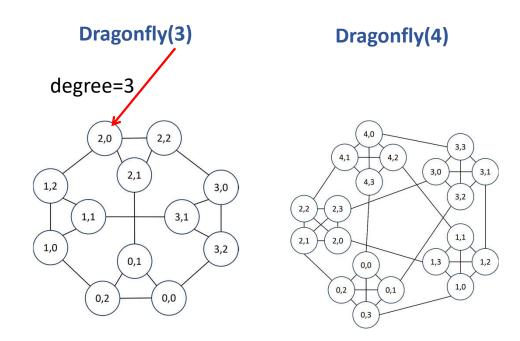




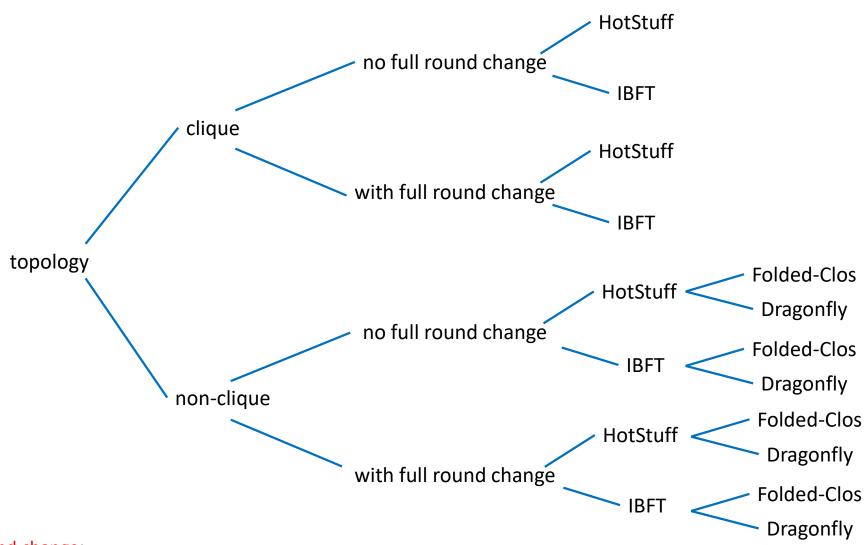
topology

"edge switch" = "switch where validators are attached"





all switches are edge switches



full round change:

all validators undergo round change for the same round (e.g.leader crashes, or validators time out).

system abstraction

application (smart contracts)

execution (VM/containers)

data (transactions)

consensus (blocks)

network (switches)

transaction throughput depends on memory pool protocol

HotStuff
Istanbul BFT (IBFT)

Clique / Folded-Clos / Dragonfly

our model

closed model:

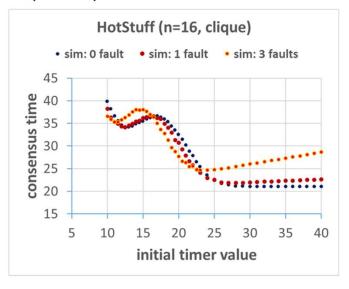
every consensus generates a new block

assume crash faults

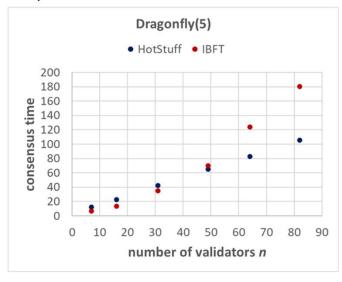
metric:

average consensus time

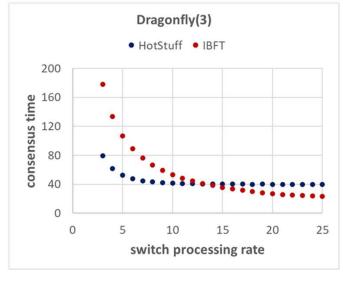
impact of parameters?



impact of scale?



impact of network?



terminology/notation:

n = #validators (participants)

$$f = [(n-1)/3]$$

 n_f = #faults $(n_f \le f)$

$$r = \frac{n_f}{n}$$
 (fraction of faults)

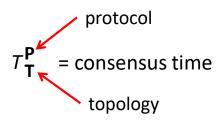
 τ_0 = initial timer value

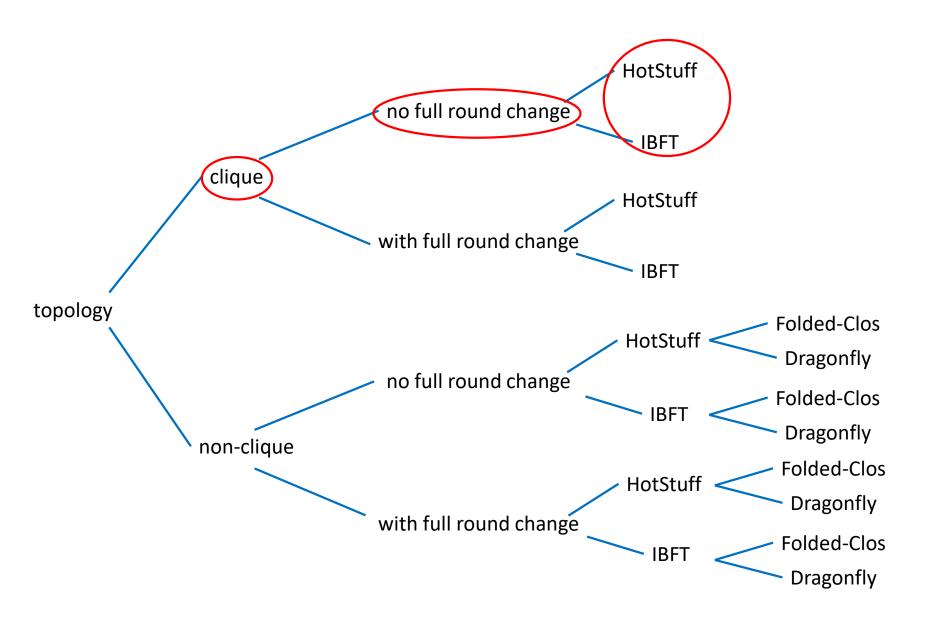
change of leader

q = Prob(full round change | nonfaulty leader)

message processing rates

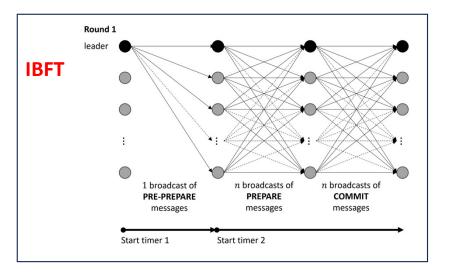


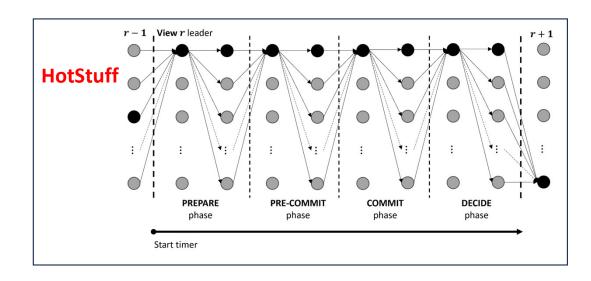




IBFT vs HotStuff (clique; no full round change)

which is faster? scalability?

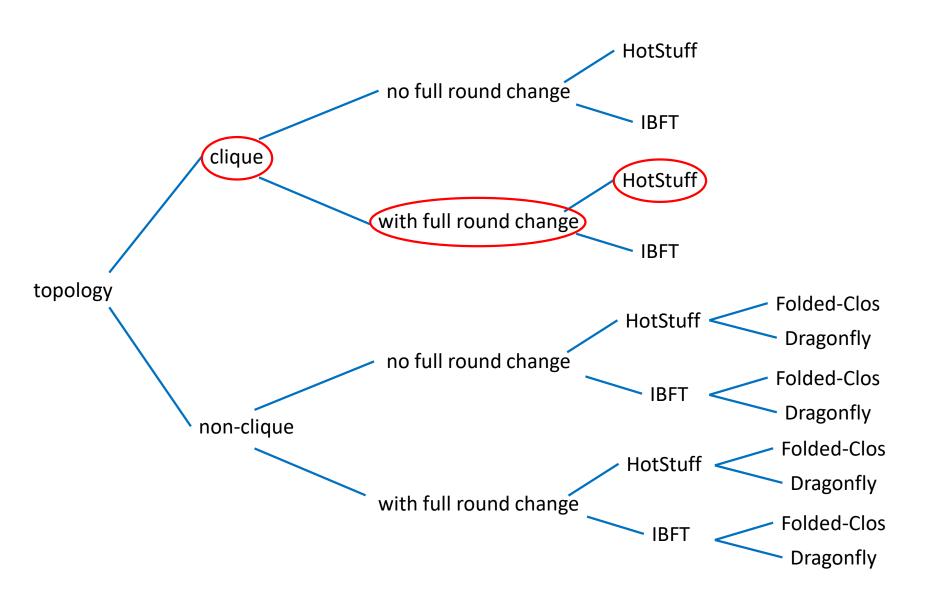




$$ET_{\mathbf{C}}^{\mathbf{I}} = \frac{2n+1}{r_{\mathbf{V}}^{\mathbf{I}}}$$

$$ET_{\mathbf{C}}^{\mathbf{H}} = \frac{4n-f+1}{r_{\mathbf{V}}^{\mathbf{H}}}$$

$$\frac{2}{r_{\mathbf{V}}^{\mathbf{H}}}$$
depends



HotStuff (clique; with full round change)

$$m^{\mathbf{H}} = 4n - 3n_f - f + 1$$
 $(n_f = \#faults)$

20

initial timer value

15

$$E(T_{\mathbf{C}}^{\mathbf{H}}) = \frac{m^{\mathbf{H}}}{r_{\mathbf{V}}^{\mathbf{H}}} + \frac{r + (1 - r)q}{1 - 2r} \tau_{0}$$

45

40

35

30

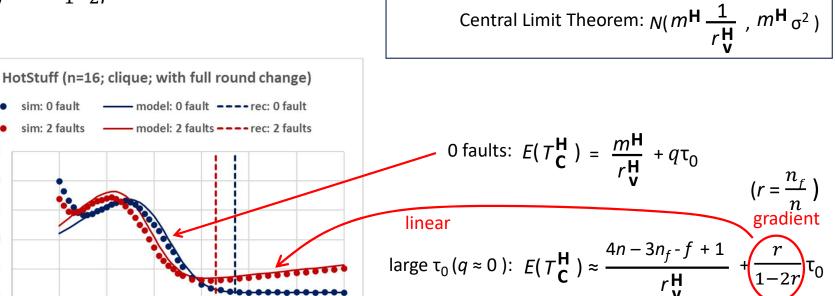
25

20

15 ^{__} 5

10

consensus time

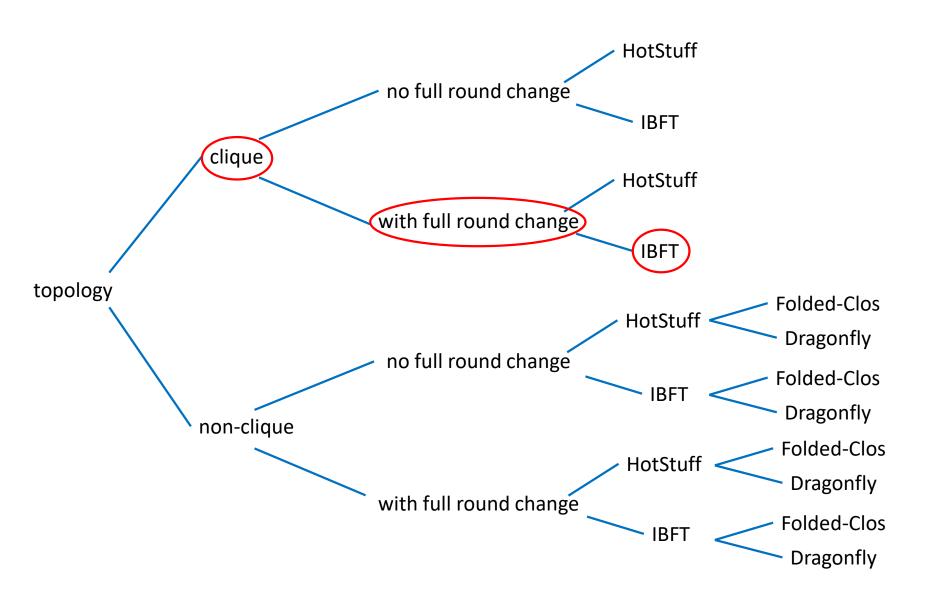


 X_i = validator processing time for *i*-th message

q = Prob(full round change | nonfaulty leader)

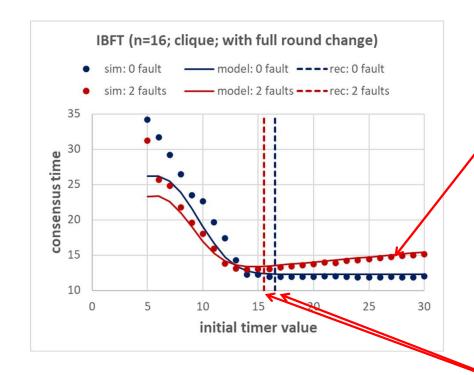
= Prob($X_1 + X_2 + ... + X_{mH} > \tau_0$)

recommended
$$\tau_0^* = mH \frac{1}{r_V^H} + 3\sqrt{mH} \sigma$$



IBFT (clique; with full round change)

$$E(T_{\mathbf{C}}^{\mathbf{I}}) = \frac{m^{\mathbf{I}}}{r_{\mathbf{V}}^{\mathbf{I}}} + \frac{r + (1 - r)q}{1 - 2r} \tau_{0} + \frac{r + (2 - r)(1 - r)q}{r_{\mathbf{V}}^{\mathbf{I}}} (n - n_{f})$$



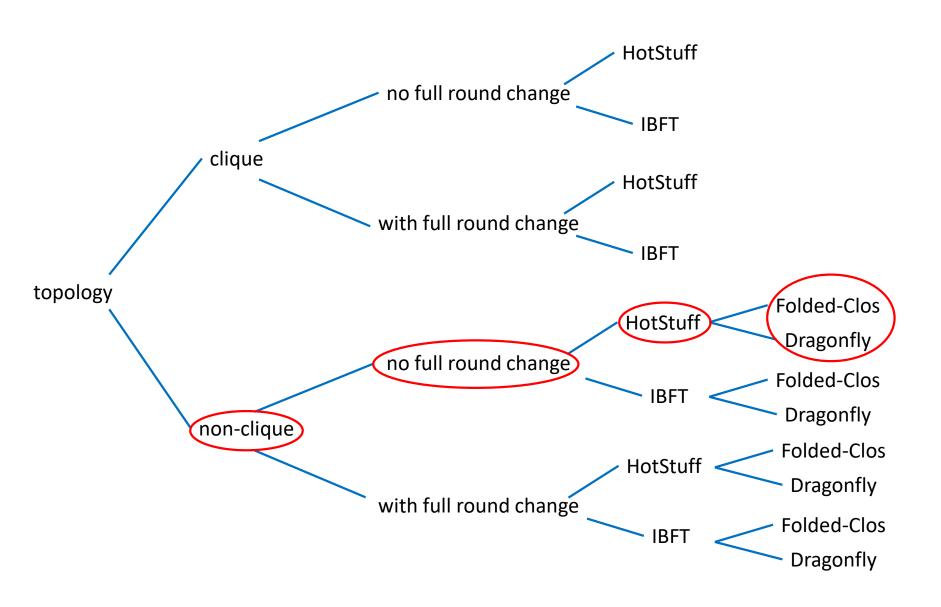
 $q = \text{Prob}(\text{full round change} \mid \text{nonfaulty leader})$ $= \text{Prob}(Y > \tau_0)$ $N(m^{\parallel} \frac{1}{r_{\text{v}}^{\parallel}}, m^{\parallel} (1 + (1 + \frac{1}{n_w}))\sigma^2)$ $n_w + n - f$ added variance from varying validator progress

large
$$\tau_0 (q \approx 0)$$
: $\frac{r}{1-2r} \tau_0$ same as HotStuff

HotStuff:
$$\tau_0^* = m^H \frac{1}{r_v^H} + 3\sqrt{m^H} \sigma$$

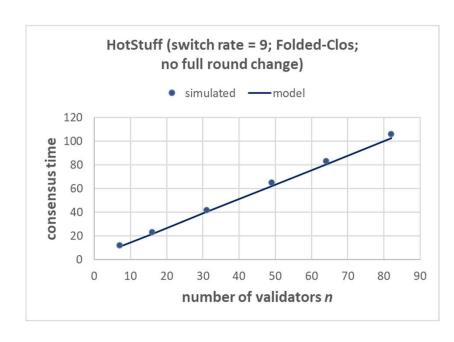
$$m^H = 4n - 3n_f f + 1$$

$$m^I = 2n - 2n_f f$$
IBFT: $\tau_0^* = m^I \frac{1}{r_v^H} + 3\sqrt{m^I(2 + \frac{1}{n_w})} \sigma$



HotStuff (nonclique; no full round change)

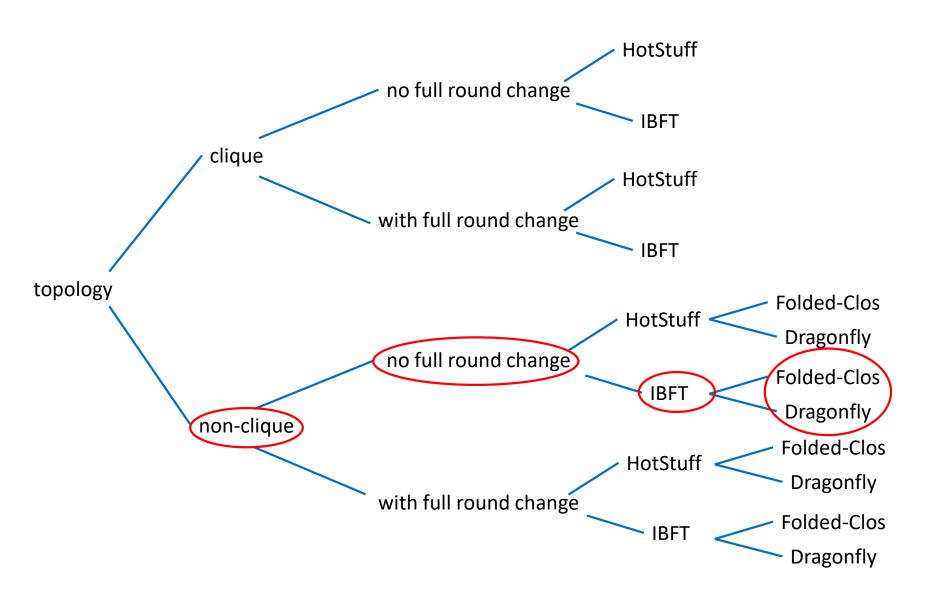
$$E(T_{\mathbf{T}}^{\mathbf{H}}) = 4\max\{\frac{n-f-1}{r_{\mathbf{V}}^{\mathbf{H}}}, \frac{n-2}{r_{\mathbf{S}}^{\mathbf{H}}}\} + 3\max\{\frac{f+2}{r_{\mathbf{V}}^{\mathbf{H}}}, \frac{n}{r_{\mathbf{S}}^{\mathbf{H}}}\} + 2(\frac{1}{r_{\mathbf{V}}^{\mathbf{H}}} + \frac{h_{\mathbf{T}}}{r_{\mathbf{S}}^{\mathbf{H}}})$$



hop count

analysis

- $E(T_{\mathbf{T}}^{\mathbf{H}})$ mostly network delay if $r_{\mathbf{s}}^{\mathbf{H}} < \frac{3}{2} r_{\mathbf{v}}^{\mathbf{H}}$ (independent of n)
- topology has negligible impact if hop count $h_T \ll n$



IBFT vs HotStuff (Dragonfly; no full round change)

suppose switch rate dominates performance:

HotStuff

$$E(T_{\mathsf{T}}^{\mathsf{H}}) = 4\max\{\frac{n-f-1}{r_{\mathsf{V}}^{\mathsf{H}}}, \frac{n-2}{r_{\mathsf{V}}^{\mathsf{H}}}\} + 3\max\{\frac{f+2}{r_{\mathsf{V}}^{\mathsf{H}}}, \frac{n}{r_{\mathsf{S}}^{\mathsf{H}}}\} + 2(\frac{1}{r_{\mathsf{V}}^{\mathsf{H}}} + \frac{h_{\mathsf{T}}}{r_{\mathsf{S}}^{\mathsf{H}}})$$

IBFT

$$E(T_{\mathbf{D}}^{\mathbf{I}}) = \max\{\frac{m_{\mathbf{V}}^{\mathbf{I}}}{r_{\mathbf{V}}^{\mathbf{I}}}, \frac{m_{\mathbf{D}}^{\mathbf{I}}}{r_{\mathbf{S}}^{\mathbf{I}}}\} \quad m_{\mathbf{D}}^{\mathbf{I}} = 2(\frac{n}{c})(2(n-1) - \frac{n}{c}(\frac{n}{c} - 1) + 2(\frac{n}{c})d(d-1)) + (n-1)$$

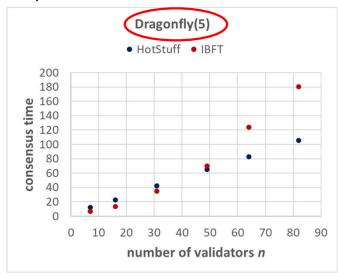
$$O(n^{2})$$

suppose we scale #switches by $c = \frac{n}{k}$, k constant

$$m_{\mathbf{D}}^{\mathbf{I}} = 2k (2(n-1) - k(k-1) + 2n \frac{d(d-1)}{d(d+1)}) + (n-1)$$

$$O(n)$$

impact of scale?



#edge-switches c = #switches

IBFT vs HotStuff (nonclique; no full round change)

case: HotStuff dominated by computation time and IBFT dominated by switching time

$$E(T_{\mathbf{D}}^{\mathbf{H}}) = 4 \max \{ \frac{n-f-1}{r_{\mathbf{V}}^{\mathbf{H}}} \} + 3 \max \{ \frac{f+2}{r_{\mathbf{V}}^{\mathbf{H}}} \} + 2 (\frac{1}{r_{\mathbf{V}}^{\mathbf{H}}}) + \frac{h_{\mathbf{D}}}{r_{\mathbf{S}}^{\mathbf{H}}} \}$$

$$= 4 \frac{n-f-1}{r_{\mathbf{V}}^{\mathbf{H}}} + 3 \frac{f+2}{r_{\mathbf{V}}^{\mathbf{H}}} + 2 \frac{1}{r_{\mathbf{V}}^{\mathbf{H}}}$$

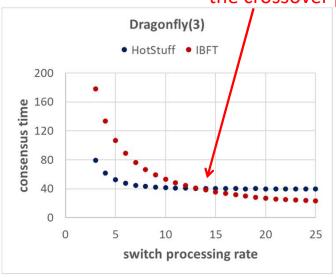
$$E(T_{\mathbf{D}}^{\mathbf{I}}) = \max\{\frac{m_{\mathbf{V}}^{\mathbf{I}}}{r_{\mathbf{V}}^{\mathbf{I}}}, \frac{m_{\mathbf{D}}^{\mathbf{I}}}{r_{\mathbf{S}}^{\mathbf{I}}}\} = \frac{m_{\mathbf{D}}^{\mathbf{I}}}{r_{\mathbf{S}}^{\mathbf{I}}}$$

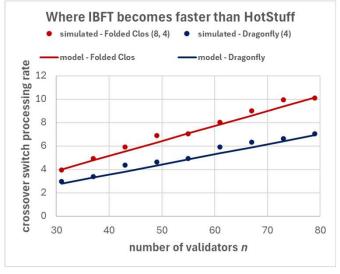
$$Arr$$
 $E(T_{\mathbf{D}}^{\mathbf{H}}) = E(T_{\mathbf{D}}^{\mathbf{I}})$ when $r_{\mathbf{S}}^{\mathbf{I}} \approx \frac{m_{\mathbf{D}}^{\mathbf{I}}}{4n - f + 4} r_{\mathbf{V}}^{\mathbf{H}}$ for Dragonfly

Similarly

$$E(T_{\mathbf{F}}^{\mathbf{H}}) = E(T_{\mathbf{F}}^{\mathbf{I}})$$
 when $r_{\mathbf{S}}^{\mathbf{I}} \approx \frac{m_{\mathbf{F}}^{\mathbf{I}}}{4n - f + 4} r_{\mathbf{V}}^{\mathbf{H}}$ for Folded-Clos

what determines the crossover point?





Folded-Clos vs Dragonfly (IBFT; no full round change)

topology has minimal impact if hopcount << n

HotStuff

$$E(T_{\mathsf{T}}^{\mathsf{H}}) = 4\max\{\frac{n-f-1}{r_{\mathsf{v}}^{\mathsf{H}}}, \frac{n-2}{r_{\mathsf{s}}^{\mathsf{H}}}\} + 3\max\{\frac{f+2}{r_{\mathsf{v}}^{\mathsf{H}}}, \frac{n}{r_{\mathsf{s}}^{\mathsf{H}}}\} + 2(\frac{1}{r_{\mathsf{v}}^{\mathsf{H}}} + \frac{h_{\mathsf{T}}}{r_{\mathsf{s}}^{\mathsf{H}}})$$

Folded-Clos and Dragonfly
have similar consensus time for HotStuff

IBFT

$$E(T_{\mathbf{T}}^{\mathbf{I}}) = \max\{\frac{m_{\mathbf{V}}^{\mathbf{I}}}{r_{\mathbf{V}}^{\mathbf{I}}}, \frac{m_{\mathbf{T}}^{\mathbf{I}}}{r_{\mathbf{S}}^{\mathbf{I}}}\}$$

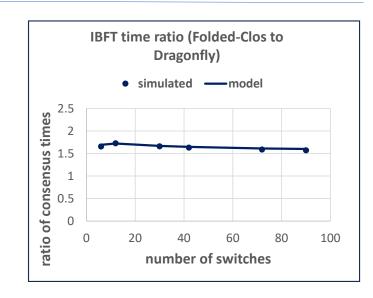
topology has minimal impact if computation time dominates

if communication time dominates (k_F , k_D = #validators per edge-switch):

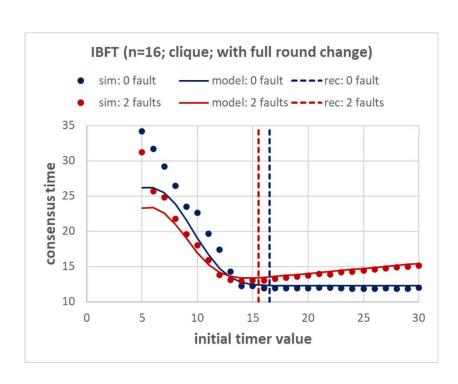
$$\frac{E(T_{\mathbf{F}}^{\mathbf{I}})}{E(T_{\mathbf{D}}^{\mathbf{I}})} = \frac{m_{\mathbf{F}}^{\mathbf{I}}}{m_{\mathbf{D}}^{\mathbf{I}}} = \frac{2k_F(2(n-1)-k_F(k_F-1))+(n-1)}{2k_D(2(n-1)-k_D(k_D-1)+2k_Dd(d-1))+(n-1)} \approx \frac{12}{8+(1/k_D)}$$

for large *n* and the same #switches

ightharpoonup for fixed k, Folded-Clos is consistently slower than Dragonfly



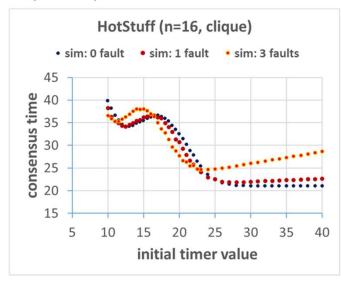
 \bullet faults (n_f) affect consensus time non-monotonically



$$E(T_{\mathbf{C}}^{\mathbf{I}}) = \frac{m!}{r_{\mathbf{V}}^{\mathbf{I}}} + \frac{r + (1 - r)q}{1 - 2r} t_{0} + \frac{r + (2 - r)(1 - r)q}{r_{\mathbf{V}}^{\mathbf{I}}} (n \cdot n_{f})$$

- \bullet faults (n_f) affect consensus time non-monotonically
- fault tolerance (τ_0) affects consensus time non-monotonically

impact of parameters?



$$E(T_{\mathbf{C}}^{\mathbf{H}}) = \frac{m^{\mathbf{H}}}{r_{\mathbf{V}}^{\mathbf{H}}} + \frac{r + (1 - r)q}{1 - 2r} \tau_{0}$$

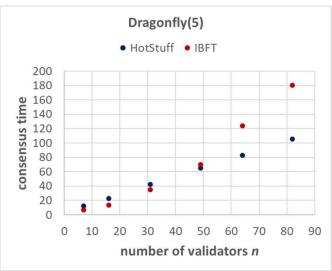
- \bullet faults (n_f) affect consensus time non-monotonically
- fault tolerance (τ_0) affects consensus time non-monotonically
- which protocol scales better?

suppose we scale #switches by $c = \frac{n}{k}$, k constant

$$m_{\mathbf{D}}^{\mathbf{I}} = 2k (2(n-1) - k(k-1) + 2n \frac{d(d-1)}{d(d+1)}) + (n-1)$$

$$O(n)$$

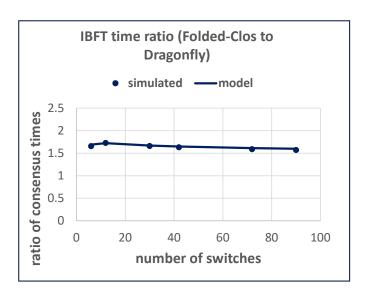
impact of scale?



- ullet faults (n_f) affect consensus time non-monotonically
- fault tolerance (τ_0) affects consensus time non-monotonically
- which protocol scales better? depends on how the network scales

- faults (n_f) affect consensus time non-monotonically
- fault tolerance (τ_0) affects consensus time non-monotonically
- which protocol scales better?
- which topology is faster?

depends on how the network scales



$$\frac{E(T_{\mathbf{F}}^{\mathbf{I}})}{E(T_{\mathbf{D}}^{\mathbf{I}})} \approx \frac{12}{8 + (1/k_D)}$$
#validators/edge-switch

- faults (n_f) affect consensus time non-monotonically
- fault tolerance (τ_0) affects consensus time non-monotonically
- which protocol scales better? depends on how the network scales
- which topology is faster for IBFT? depends on #validators/edge-switch

Y.C. Tay

(keynote for IWSF&SHIFT 2025, Sao Paolo, Brazil)