

1 Object Oriented Consensus

2 **Yehuda Afek**

3 Department of Computer Science, Tel-Aviv University, Israel
4 afek@tau.ac.il

5 **James Aspnes**

6 Department of Computer Science, Yale University, USA
7 james.aspnes@gmail.com

8 **Edo Cohen**

9 Department of Computer Science, Tel-Aviv University, Israel
10 edocohen@tau.ac.il

11 **Danny Vainstein**

12 Department of Computer Science, Tel-Aviv University, Israel
13 dannyvainstein@gmail.com

14 — Abstract —

15 Our work focuses on the problem of decomposing consensus algorithms into a common frame-
16 work composed of simple building blocks. We show that earlier decomposition strategies fall
17 short when applied to some well known algorithms and present a new framework in order to
18 tackle the problem. First we use Aspnes' framework [2] composed of *adopt-commit* [5] and *con-*
19 *ciliator* [2] objects in order to decompose the well known Phase-King Byzantine algorithm [4].
20 We then consider two other well-known algorithms and argue that this framework is insufficient
21 in these (and other) cases and offer a new framework. The framework works in rounds where
22 each consists of two steps. The first step involves an object which detects agreement and the
23 second involves an object that aims at achieving consensus. We denote our newly defined objects
24 as *vacillate-adopt-commit* and *reconciliator*. We demonstrate our decomposition on two
25 well known algorithms. Namely, Ben-Or's Randomized algorithm [3] and the Raft algorithm [6].

26 **2012 ACM Subject Classification** Theory of computation → Models of computation → Concur-
27 rency → Distributed computing models

28 **Keywords and phrases** distributed algorithms, wait-free, consensus, message-passing

29 **Digital Object Identifier** 10.4230/LIPIcs.CVIT.2016.23

30 **1 Introduction**

31 The consensus problem introduced by Lamport, Pease and Shostak [7] resides at the heart of
32 many distributed algorithms such as leader election, database transaction handling, resource
33 allocation, ensuring storage replicas are mutually consistent, block chain technology and
34 many more.

35 In the consensus protocol between n processors, each with an input value, processors
36 agree on a single common output which was the input to one of them. While consensus
37 is trivial in a non-faulty synchronous environment, it is often more difficult in practice as
38 most distributed networks are asynchronous and must be resilient to faults and various miss
39 behaviors.

40 The philosophy of software engineering asserts that decomposing complex systems into
41 simple building blocks is a good thing. One reason for this being that by analyzing simpler
42 objects we may deduce observations on more complex systems.



© Y. Afek, J. Aspnes, E. Cohen, D. Vainstein;
licensed under Creative Commons License CC-BY
42nd Conference on Very Important Topics (CVIT 2016).

Editors: John Q. Open and Joan R. Access; Article No. 23; pp. 23:1–23:16



Leibniz International Proceedings in Informatics

Schloss Dagstuhl – Leibniz-Zentrum für Informatik, Dagstuhl Publishing, Germany

43 One attempt to decompose consensus was provided by Gafni [5]. Gafni proposed *adopt-*
 44 *commit*, an object fulfilling weaker guarantees than consensus, as a building block of consensus.
 45 Aspnes further provided a detailed decomposition [2] of consensus into *adopt-commit* and
 46 a complementary *conciliator* object which together form a generic framework describing
 47 consensus. In this work we extend these decompositions, into *vacillate-adopt-commit* and
 48 *reconciliator* objects, which describe the core of many existing consensus algorithms in which
 49 it is not clear how to breakdown into the previous frameworks. Our main thesis in the paper
 50 is that many known consensus algorithms fall into a similar pattern of a repetitive two-fold
 51 process in which the first step evaluates the current status and how close it is to a consensus
 52 and a second step which brings closer to a consensus by taking some action. Together with
 53 previous work, we provide a more solid interpretation of the mechanisms composing well
 54 known consensus algorithms.

55 Our paper is structured as follows. In Section 3, the consensus template is presented. In
 56 section 4 we first tackle the decomposition of the well known phase-king byzantine algorithm
 57 [4] and show how it naturally decomposes Aspnes' framework of *adopt-commit* and *conciliator*
 58 objects. Thereafter we tackle Ben-Or's randomized algorithm [3] and the raft algorithm [6].
 59 We observe that Aspnes' former framework fails to capture these algorithms naturally and
 60 provide our alternate decomposition using *vacillate-adopt-commit* and *reconciliator* objects.

61 In section 5 we explore the relation between *vacillate-adopt-commit* and *adopt-commit*,
 62 showing the latter is slightly weaker. We conclude with final remarks in Section 6.

63 2 Preliminaries

64 We consider a network consisting of n processors, $\{p_1, \dots, p_n\}$, each processor p begins with
 65 an input value $p.init$. The goal of **consensus** protocol is to agree on a value, u , for all
 66 processors while satisfying the following three conditions:

- 67 ■ **Validity** - The value u must be the initial value of one of the processors involved in the
 68 protocol, i.e., $\exists i, s.t. p_i.init == u$.
- 69 ■ **Termination** - Each processor decides after taking a finite number of steps.
- 70 ■ **Agreement** - The value agreed upon, u , is the same for all processors.

71 In the context of distributed computation it is useful to refer to an implementation of a
 72 protocol as an **object**. We follow these notations in this paper, that is, a processor which is
 73 part of a network attempting to achieve consensus invokes an object, $Consensus(v)$, with its
 74 initial value $v \leftarrow p.init$ and expects to receive a value u such that the properties above are
 75 fulfilled.

76 An **adopt-commit** object may be viewed as a weaker version of consensus. In **adopt-**
 77 **commit** the returned bit is accompanied by a confidence level. The **agreement** property
 78 is waived and two more refined properties are defined instead. Formally, **adopt-commit**
 79 receives as an argument a value v and returns a pair $(confidence, u)$ where $confidence \in$
 80 $\{adopt, commit\}$ fulfilling the following properties:

- 81 ■ **Validity, Termination** - Same as above.
- 82 ■ **Coherence** - If some processor receives $(commit, u)$ then all other processors receive the
 83 same value u (with confidence either *adopt* or *commit*).
- 84 ■ **Convergence** - If all processors invoke **adopt-commit**(v) with the same value v , then
 85 all processors receive $(commit, v)$.

86 The newly introduced object **vacillate-adopt-commit** adds an additional confidence
 87 level with respect to **adopt-commit**. The object expects a value v from every processor and

88 returns a value u and one of three confidence levels, *vacillate*, *adopt* or *commit*. The
 89 following guarantees are required of the object:

- 90 ■ **Validity, Termination, Convergence** - Same as above.
- 91 ■ **Coherence over adopt & commit** - If any processor received $(commit, u)$, then every
 92 other processor receives either $(commit, u)$ or $(adopt, u)$ (this guarantee corresponds to
 93 the adopt-commit's coherence).
- 94 ■ **Coherence over vacillate & adopt** - If no processor received *commit* and some
 95 processor received $(adopt, u)$, then every other processor receives either $(adopt, u)$ or
 96 $(vacillate, *)$, where $*$ may be any value (subject to other constraints like validity).

97 Our main thesis in this paper is that consensus algorithms are essentially a sequence of
 98 repetitive rounds of invoking an object which revises the current status of the processors and
 99 how close they are to an agreement followed by an action which brings the processors closer
 100 to an agreement. The objects *adopt-commit* and *vacillate-adopt-commit* are building blocks
 101 suitable to fulfill the role of the agreement detector. In order to provide all the building
 102 blocks necessary to compose consensus we must define objects which perform some action
 103 that possibly brings closer towards a consensus. In his framework [2], Aspnes defines such an
 104 object, *Conciliator*, which is required to satisfy the following properties:

- 105 ■ **Validity, Termination** - Same as above.
- 106 ■ **Probabilistic Agreement** - With probability greater than 0, the returned value is the
 107 same to all processors.

108 Another object we define is the *Reconciliator*. The *reconciliator* is required to satisfy
 109 the following properties:

- 110 ■ **Termination** - Same as above.
- 111 ■ **Weak Agreement** - With probability 1, at some point all invoking processors receive
 112 the same value such that this value corresponds to the *adopt* values of the current round.¹

113 Despite the name similarities, these objects differ profoundly, the conciliator is essentially
 114 a probabilistic consensus while the reconciliator is weaker in the sense that it may be invoked
 115 by a subset of the network.

116 Throughout the paper, we abbreviate the objects *vacillate-adopt-commit* as **VAC** and
 117 *adopt-commit* as **AC**. Furthermore, the returned values are denoted by the first letter of the
 118 confidence levels, i.e. **V, A, C**.

119 **3 The Generic Form of Consensus**

120 The idea behind this paper is that many well known consensus algorithms have the same
 121 basic structure consisting of two objects, *VAC* that checks whether consensus has been
 122 reached or not and *reconciliator* that shakes up the preferences of the processors in case of
 123 a stalemate. We do not claim that this structure is necessarily better than using *AC* and
 124 *conciliators*, but that it more accurately reflects the existing structure of algorithms in the
 125 field.

126 Informally, the generic consensus algorithms work in rounds. In each round, first the
 127 *VAC* is invoked to observe the system state and inform us as to whether consensus has been
 128 reached. *VAC* returns to each processor one of three possible outputs: (1) $(commit, v)$

¹ if no such values were given, the returned value should be one of the inputted values

129 which indicates that the system has reached an agreement on value v , (2) (*adopt*, v) which
 130 indicates that it is possible that some processors in the system have agreed on the value v ,
 131 and (3) (*vacillate*, v) indicating that the system is in an indecisive state.

132 If a processor receives (*commit*, v), it is guaranteed that no other processor receives a
 133 *vacillate* value and all outputs return with the same value, v . A processor receiving (*adopt*,
 134 v) is guaranteed that any other processor either received a *vacillate* value or received the
 135 same preference that the earlier processor has. Finally, if a processor receives (*vacillate*,
 136 v), the only guarantee it has is that no other processor received a *commit* value.

137 We note the key difference between the *VAC* and *AC* objects. An adopt-commit object
 138 always returns a new value to be adopted by a process, but this is not consistent with the
 139 structure of many consensus protocols in the literature. Adding a third option, i.e., *vacillate*,
 140 accounts for situations where the algorithm does not force a process to update its preference.
 141 Furthermore *vacillate* gives the receiving processor the information that a consensus has not
 142 been reached. This type of information is not available using the *AC* infrastructure, but as
 143 we will see is available in existing algorithms in the literature.

144 The question is how termination of the consensus can be guaranteed if the collection of
 145 preferences is balanced and the *VAC* continually returns *vacillate*. For that purpose, the
 146 *reconciliator* is used to give each *vacillating* processor a new preference with a guarantee to
 147 provide a deciding set of preferences with some probability. That is, whenever a processor
 148 receives a *vacillate* value from the *VAC* object the distributed *reconciliator* provides each
 149 processor with an alternate preference which guarantees that eventually enough processors
 150 will get the same preference leading to *VAC* eventually observing agreement. Pseudocode for
 151 the consensus template is given in Algorithm 1.

```

Consensus( $v$ )
   $m \leftarrow 0$ 
  INIT()
  while true do
     $m \leftarrow m + 1$ 
     $(X, \sigma) \leftarrow VAC(v, m)$ 
    switch  $X$  do
      case vacillate:
         $v \leftarrow Reconciliator(X, \sigma, m)$ 
      case adopt:  $v \leftarrow \sigma$ 
      case commit:  $v \leftarrow \sigma$  and
        decide  $\sigma$ 
    endsw
  end

```

Algorithm 1: Consensus Template

```

Consensus( $v$ )
   $m \leftarrow 0$ 
  INIT()
  while true do
     $m \leftarrow m + 1$ 
     $(X, \sigma) \leftarrow AC(v, m)$ 
    switch  $X$  do
      case adopt:
         $v \leftarrow Conciliator(X, \sigma, m)$ 
      case commit:  $v \leftarrow \sigma$  and
        decide  $\sigma$ 
    endsw
  end

```

Algorithm 2: Consensus Template using AC and Conciliator

152 Note that *INIT* is a void function unless stated otherwise. Furthermore, note that the
 153 operation, *decide* σ , is followed by a halt operation, that is, the processor will decide upon
 154 its value and return. The argument m is the phase of the consensus process.

155 Next we prove that the template indeed achieves consensus, using the *VAC* and *reconciliator*
 156 properties.

157 ► **Lemma 1.** *Algorithm 1 is a correct consensus algorithm.*

158 **Proof. ■ Agreement:** If some processor decided on a value v we are guaranteed that it
 159 received ($commit, v$) during that round. Thus by the VAC's coherence over adopt commit
 160 guarantee we are ensured that all processors complete the round with the same value.
 161 Thus by VAC's convergence all remaining processors will decide on the same value in the
 162 next round.

163 ■ **Validity:** Follows from the reconciliator's guarantees and the fact that all inputs to the
 164 VAC and reconciliator objects are clearly valid inputs.

165 ■ **Termination:** Follows from the reconciliator's guarantees.

166

167 In addition to algorithm 1, we provide here a concrete algorithm for consensus using the
 168 AC and conciliator objects. The algorithm correctness is proved similarly to 3.

169 4 Consensus Decomposition

170 We first demonstrate how the phase-king algorithm naturally decomposes using Aspnes'
 171 framework [2] (i.e. adopt-commit rather than vacillate-adopt-commit and conciliator rather
 172 than reconciliator). Following that we decompose ben-or's algorithm and the raft algorithm
 173 using our framework under the statement that Aspnes' earlier framework is not sufficient.

174 4.1 Phase-King Algorithm

175 Here we show how the Phase-King consensus protocol of Berman, Garay and Perry [4]
 176 fits the framework given in [2]. Throughout this section we assume a message passing
 177 synchronous model and t byzantine processors such that $3t < n$. Note that in contrast
 178 to the original consensus template, every algorithm continues to participate in the overall
 179 consensus template even after deciding upon a value, as in the original Phase-King algorithm.
 180 Algorithms 3 and 4 are the AC's and *conciliator*'s implementations and lemmas 2 and 3
 181 prove the implementations correctness, i.e., ensuring the objects' guarantees.

182 ► **Lemma 2.** *Algorithm 3 is a correct adopt-commit implementation.*

183 **Proof.** The proof is similar to the Phase-King correctness proof given in [4].

184 ■ **Validity:** If all inputs are the same value $v \in \{0, 1\}$ then at the end of the first exchange
 185 of the first round all $C(v)$ values are at least $n - t$ for all correct processors (since there
 186 are $n - t$ non-byzantine processors) and $C(u) \leq t < n - t$ for $u \neq v$ (since $t < n/3$).
 187 Therefore v will remain the chosen value for all non-byzantine processors after exchange 1.
 188 The same holds after exchange 2. Thus all correct processors will enter the *if* statement
 189 and receive a value of ($commit, v$) guaranteeing validity.

190 ■ **Convergence:** If all inputs are the same value v then but what we have shown earlier
 191 all processors will receive a value of ($commit, v$) as needed.

192 ■ **Termination:** Since we are in the synchronous setting clearly the protocol terminates
 193 within a finite amount of rounds.

194 ■ **Coherence over commit and adopt:** We first observe that after exchange 1 exists
 195 some value $v \in \{0, 1\}$ such that all correct processors' values are either 2 or v (this
 196 follows since $t < n/3$ meaning that otherwise we would have some correct processor that
 197 broadcasted different values to different processors). Next we observe that in fact this
 198 property holds through exchange 2 as well (again, since otherwise we would have a correct
 199 processor that broadcasted different values to different processors). Thus, if 2 *commit*
 200 messages were received (meaning their values are $v \neq 2$) then their values must coincide.

```

AC(v,m)
  broadcast ⟨v⟩          // (*
  Exchange 1 *)
  v ← 2
  for k=0 to 1 do
    C(k) ← # received k's
    if C(k) ≥ n - t then
      v ← k
    end
  end
  broadcast ⟨v⟩          // (* Exchange
  2 *)
  for k=2 downto 0 do
    D(k) ← # received k's
    if D(k) > t then
      v ← k
    end
  end
  if (v ≠ 2 and D(v) ≥ n - t) then
    return (commit, v)
  else
    return (adopt, v)
  end

```

Algorithm 3: Phase-King's *adopt-commit* implementation

```

Conciliator(X, σ, m)
  if id = m then
    broadcast ⟨MIN(1, v)⟩
    σm ← received message from
    processor m
  return (adopt, σm)

```

Algorithm 4: Phase-King's *conciliator* implementation

201

202 ► **Lemma 3.** *Algorithm 4 is a correct conciliator implementation.*

203 **Proof.** ■ **Validity** Follows since the phase king's inputted value is σ_m .

204 ■ **Termination** Clearly guaranteed.

205 ■ **Coherence** Only *adopt* is returned therefore coherence holds.

206 ■ **Probabilistic Agreement** Since our setting are not probabilistic but rather deterministic we will show eventual agreement, i.e., eventually the conciliator will cause agreement. Consider round m such that processor m is non-byzantine. If during this round all returned values from the adopt-commit object were *adopt* values then following this *conciliator* round all object adopt the same value from the phase-king. Otherwise some processor received at least $n - t$ values that are different than 2 during exchange 2. Since $t < n/3$ this means that $> t$ such values were broadcasted during that exchange. Since as we have seen there is only one value $\neq 2$ that is adopted during exchange 2 and since there are only t byzantine processors, p_m will have adopted that same value and therefore would have broadcasted it through the *conciliator* as needed.

216

217 4.2 Ben-Or's Algorithm

218 In this section we consider Ben-Or's algorithm [3]. Throughout this section the settings are
 219 asynchronous, message-passing model and the number of tolerated crash failures, t is strictly

220 smaller than $n/2$.

221 After some consideration regarding the decomposition of this algorithm we came to
 222 the conclusion that the use of a single adopt-commit followed by a conciliator object is
 223 insufficient. This is mainly due to the fact that in Ben-Or's algorithm there are 3 unique
 224 types of processor per round - processors that received more than t ratify message, processors
 225 that received atleast 1 ratify message but less than t and processors that received no ratify
 226 messages. Each type of processor has some guarantee about the state of the network (e.g., if
 227 received more than t ratify messages one is ensured the network has achieved consensus).
 228 More regarding why adopt-commit is insufficient given in section 5.

229 Algorithms 6 and 5 are the *VAC*'s and *reconciliator*'s implementations, Lemmas 4 and 5
 230 prove the implementations correctness, i.e., that they uphold the objects' guarantees.

<pre> VAC(v, m) send $\langle 1, v \rangle$ to all wait to receive $n - t$ $\langle 1, * \rangle$ messages if received more than $n/2$ $\langle 1, v \rangle$ messages then send $\langle 2, v, ratify \rangle$ to all else send $\langle 2, ? \rangle$ to all end wait to receive $n - t$ $\langle 2, * \rangle$ messages if received more than t $\langle 2, v, ratify \rangle$ messages then return ($commit, v$) else if received a $\langle 2, v, ratify \rangle$ message then return ($adopt, v$) else return ($vacillate, v$) end </pre>	<pre> Reconciliator(X, σ, m) return $CoinFlip()$ </pre>
---	--

Algorithm 6: Ben-Or's *reconciliator* implementation

Algorithm 5: Ben-Or's *vacillate-adopt-commit* implementation

231 ► **Lemma 4.** *Algorithm 6 is a correct reconciliator implementation.*

232 **Proof.** Since any value has a non-zero probability of being outputted, the reconciliator's
 233 guarantee clearly follows. ◀

234 ► **Lemma 5.** *Algorithm 5 is a correct vacillate-adopt-commit implementation.*

235 **Proof.** The proof is similar to the Ben-Or algorithm correctness proof found in the survey of
 236 Aspnes [1].

237 ■ **Validity:** Follows since all messages sent and received hold inputted values.

238 ■ **Termination:** Since the number of crash failures is less than half, all processors will
 239 terminate in a finite amount of time.

240 ■ **Convergence:** Assume all processors start out with the same value v . Since the number
 241 of crash failures is less than half all live processors will send a ratify message with the
 242 same value v . Thus the if statement in line 10 will clearly be satisfied and all processors
 243 will receive a value of ($commit, v$).

- 244 ■ **Coherence over adopt and commit:** If some processor received more than t (*ratify*, v)
 245 messages then by the definition of t at least one live processor broadcast a (*ratify*, v)
 246 message. Thus by line 9 we are guaranteed that every processor received at least one
 247 ratify message. Finally the condition in line 4 we are guaranteed that if 2 ratify messages
 248 are sent out then they have the same value. Thus all processors received at least one
 249 ratify message and they all received the same value.
- 250 ■ **Coherence over vacillate and adopt:** Since the condition in line 4 insures that ratify
 251 messages hold the same value, v if some processor received an (*adopt*, v) message then
 252 clearly all other *adopt*-receiving processors received the same value v .

253

254 4.3 Raft

255 In algorithm 7 we use the Raft algorithm [6] to achieve consensus. The raft algorithm is
 256 designed for producing a consistent log among distributed systems. Every processor maintains
 257 an indexed log of commands which they update continuously. Occasionally the processors
 258 apply the commands from their logs to their state machine. The commands applied are
 259 always in order and always continue from the last command applied.

260 Here we describe how the raft algorithm may be used in order to achieve consensus. We
 261 use the raft algorithm with a single command (i.e., the logs will consist of a single type of
 262 command). The single command used is decide-and-stop-applying-to-state-machine which
 263 we denote as **D&S(v)**. This command tells the state machine to decide on the value v and
 264 stop applying any further commands thereafter (i.e., not to switch its decision).

265 In the raft algorithm every processor updates an indexed log continuously and occasionally
 266 applies the commands given in the log to its state machine. The commands applied are
 267 always in order and always continue from the last command applied. Therefore, once a
 268 processor decides to update its state machine it will apply the first command in the log, and
 269 decide on that given value. This results in the processor deciding upon the first value it sees
 270 in its log.

271 The raft algorithm works as follows; there are 3 states a processor may be in, follower,
 272 candidate and leader. Every processor aspires to become a leader and once becoming a leader
 273 it tries to have the system decide upon its value. All processors start out as followers and
 274 employ a timer. The timer is reset every time the processor receives a message from a fellow
 275 processor (with the caveat of terms which will be explained next).

276 Once a timer runs out the processor converts to candidate and tries to gain enough
 277 votes to become leader. Meaning that once becoming a candidate the processor broadcasts
 278 *RequestVote* messages and if it achieves a majority of acks it converts to leader.

279 Once leader, the processor tries to have all other processors append its value to their
 280 log by sending them a **D&S(v)** command. Initially the command is sent out as a tentative
 281 command. If the leader receives a majority of acks for this append message (denoted by
 282 *AppendEntries*) it commits to the command and broadcasts the fact that this command
 283 should now be committed to.

284 Since raft is log based the processors employ a commit-index mechanism in order to
 285 achieve the formerly described attributes. An *AppendEntries* message includes a commit
 286 index. Meaning that the receiving processor appends the entries it received from the message
 287 to its log, however it does not yet apply these commands to its state machine as they may be
 288 altered. The processor then looks at the commit index and only then applies all commands
 289 up to and including that index in its log. This results in two types of *AppendEntries* messages;
 290 the first does not change the commit index of the receiving processor but rather appends

291 commands the processor's log. The second does aim at appending entries to logs but rather
292 to update the receiving processor's commit index. Since the second type is only sent out once
293 a majority of processors acknowledge the first type, we are ensured that once a command is
294 committed to it will not change (denoted as the **state machine safety** property which will
295 be formally defined later in the section).

296 The main idea in the algorithm is that in order to achieve consensus we use a 2 step
297 mechanism (not unlike Ben-Or's algorithm). Each processor first tries to gain a majority
298 which would result in leadership. Once leader it then tries to push its a value to decide upon.
299 Only once achieving a majority of acks to that operation, it commits to that log entry and
300 notifies everyone else. We note that both 'wait' operations do not hinder the algorithm's
301 termination since in the background all non-leaders have a randomized timer which has the
302 soul purpose of shaking the protocol out of a stalemate. Once a non-leader's timer runs out
303 it begins the consensus algorithm described in algorithm 7.

304 The algorithm as described would work just fine in a system without failures, however
305 this is never the case. In real world scenarios processors may fail unexpectedly and messages
306 may be delayed or even lost. In order to maintain log consistency and ensure termination
307 even under these conditions the raft algorithm introduces the notion of terms.

308 Terms are defined such that leaders are leaders only of a specific (and all lower) term.
309 Once a leader encounters a higher term it immediately reverts to follower and updates its
310 term. Furthermore once a processor converts to candidate and tries to become leader, it
311 increases its term in order to do so. This ensures us that even though some processors may
312 fail ultimately consensus will still be obtained. We note that every processors log consists of
313 indexed entries (indexed continuously from 1, i.e., 1,2,3,...) such that each entry consists of
314 a command and the term in which the command was received.

315 It may be the case that once a leader is elected it immediately crashes (or it is somehow
316 cut off from a majority of the network). Therefore a different processor will become leader
317 before the earlier one had the chance to alter the processor's logs. This may happen over and
318 over causing a cycle of leaders without any alterations in any of the logs. This would hinder
319 the termination property of our consensus protocol. Therefore the following assumption is
320 made (note that this assumption is made in the original raft paper as well).

321 We make the assumption that the broadcast time (time it takes to convey a message)
322 is much smaller than randomized timer which is in turn much smaller than the average
323 time between failures of a single machine. This constraint is required in order to maintain
324 a leader which in turn results in consensus termination. We refer to this property as the
325 *timing property*.

326 We again note the similarity between Ben-Or's algorithm and the algorithm considered
327 here. Both algorithms use a two step mechanism - the first step alerting some processor (the
328 leader in our case) that consensus has been achieved (i.e., all processor's chose the same
329 leader) and the second step conveying that information to all other processors.

330 Next we turn to formally define the algorithm. Figures 1 and 2 describe the protocol's
331 inner state variables and types of messages used. Note that arrays NextIndex and MatchIndex
332 only apply once a processor is in leader state (and last only for the duration of that term).
333 Furthermore they are reinitialized every time the processor converts to leader. Both these
334 arrays are introduced in order to maintain consistency over the processors' logs. When
335 an AppendEntry message is received the processor may reject it (and return false) if the
336 senders log does not agree with its own (up to a certain degree which will be described next),
337 therefore these 2 arrays help the leader know how far back in its log it has to send to each
338 processor in order to ensure a positive ack.

23:10 Object Oriented Consensus

339 Algorithms 7, 8 and 9 describe the algorithm of a processor that manages to convey its
340 value (algorithm 7) and the responses to the different messages (algorithms 8 and 9). Note
341 that variables marked with * represent the variables of the processor being described.

342 We also note the behaviour of processors which receive *AppendEntries* messages (shown
343 in algorithm 9). Once a processor wakes up from a crash it contains an outdated log. Thus
344 once it receives an *AppendEntries* message it may be the case where the processor will have a
345 large chunk of commands missing in its log. In order to prevent such situations the receiving
346 processor has the option of rejecting the message (by returning a false ack). Once a leader
347 receives such an ack it then uses its *MatchIndex* and *NextIndex* values in order to go back
348 in its log to the place where the receiving processor crashed. It then retries sending the
349 *AppendEntries*, however this time the message will include all entries missing in the receiving
350 processor's log (rather than just the last log entry). Therefore, the default leader behaviour
351 would be to only send its last log entry and if that is not enough it continuously retries with
352 an earlier log.

RequestVote[term, candidateId, lastLogIndex, lastLogTerm], where lastLogIndex is
the index of the processor's last log and lastLogTerm is the term of that log index.

ack_RequestVote[term, voteGranted], where voteGranted is a boolean variable.

AppendEntries[term, leaderId, prevLogIndex, prevLogTerm, D&S(v),
leaderCommit], where prevLogIndex is the index of the log preceding the D&S(v)
command and prevLogTerm is its term.

ack_AppendEntries[term, success], where success is a boolean variable.

Figure 1: Raft Consensus Messages

CurrentTerm.

VotedFor - candidateId voted for in current term.

Log[] - indexed list of commands and terms during which they were received.

CommitIndex - log index stating that all commands up and through that index are
to be applied to the state machine.

LastApplied - log index of last command applied to state machine.

State - one of follower, candidate or leader.

NextIndex[] (variable applies only while leader) - array of length n (number of
processors). Each element is the index of the next log entry to send to that processor.
Initialized after election to the leader's last log entry + 1.

MatchIndex[](variable applies only while leader) - array of length n . Each element is
the index of highest log entry known to be replicated on that server. Initialized to 0.

Figure 2: Raft Consensus Inner State Variables

353 We first state a few properties of the Raft algorithm which were stated in [6]. These were
354 also proven in the same paper and therefore, due to space constraints we omit their proofs.

355 ■ **Leader Completeness:** if a log entry is committed in a given term, then that entry
356 will be present in the logs of the leaders for all higher-numbered terms.

357 ■ **State Machine Safety:** if a server has applied a log entry at a given index to its state
358 machine, no other server will ever apply a different log entry for the same index.

359 ■ **Log Matching:** if two logs contain an entry with the same index and term, then the
360 logs are identical in all entries up through the given index.

361 Next, we turn to prove that the algorithm guarantees consensus.

```

Consensus( $v$ )
  state  $\leftarrow$  candidate
   $v^* \leftarrow \text{log}[\text{lastLogIndex}^*]$ 
  Broadcast RequestVote[ $t^*, id^*, \text{lastLogIndex}^*, \text{lastLogTerm}^*$ ]
  Wait to receive  $> n/2$ 
  ack_RequestVote( $t = t^*, \text{granted} = \text{true}$ )
  state  $\leftarrow$  leader
   $v^* \leftarrow \text{log}[\text{lastLogIndex}^*]$ 
  Broadcast AppendEntries[ $t^*, id^*, \text{prevLogIndex}^*, \text{prevLogTerm}^*, D\&S(v^*), \text{commitIndex}^*$ ]
  Wait to receive  $> n/2$ 
  ack_AppendEntries( $t = t^*, \text{success} = \text{true}$ ) // commitIndex
  is therefore increased, see leader responses
  Broadcast AppendEntries[ $t^*, id^*, \text{prevLogIndex}^*, \text{prevLogTerm}^*, -, \text{commitIndex}^*$ ]

```

Algorithm 7: Raft Consensus Protocol

```

if received acks_AppendEntries( $t, \text{false}$ )
  from  $p_i$  then
  | if  $t > t^*$  then
  | | state  $\leftarrow$  candidate, increase term
  | else
  | | decrement NextIndex[ $i$ ], resend
  | | AppendEntries.
if received acks_AppendEntries( $t, \text{true}$ )
  from  $p_i$  then
  | update NextIndex[ $i$ ] and
  | MatchIndex[ $i$ ]
  | if exists  $N$  s.t.  $N > \text{commitIndex}^*$ ,
  |  $\text{maj of MatchIndex}[i] \geq N$ ,
  |  $\text{log}[N].\text{term} = t^*$  then
  | |  $\text{commitIndex}^* \leftarrow N$ 

```

Algorithm 8: Raft Consensus Leader Responses

```

if received RequestVote[ $t, id, \text{lastLogIndex}, \text{lastLogTerm}$ ] then
  | if ( $t < t^*$ ) || ( $t = t^* \&\& \text{votedFor} \neq \text{null}$ ) then
  | | return ( $t^*, \text{false}$ )
  | if  $\text{VotedFor} == \text{null} \&\& \text{log matches requestor's}$  then
  | | return ( $t^*, \text{true}$ ), update  $t^*$ 
if received AppendEntries[ $t, -, \text{prevLogIndex}, \text{prevLogTerm}, D\&S(v), \text{commitIndex}$ ] then // tentatively
  log  $D\&S(v)$ 
  | if  $t < t^*$  then
  | | return ( $t^*, \text{false}$ )
  | else if log does not match requestor's log at  $\text{prevLogIndex}$  then
  | | return ( $t^*, \text{false}$ )
  | else
  | | append new entries, delete conflicting ones, if deleted delete all entries that follow as well
  | |  $\text{commitIndex}^* \leftarrow \min(\text{leaderCommit}, \text{index of last new entry})$ 
if received AppendEntries[ $t, -, \text{prevLogIndex}, \text{prevLogTerm}, -, \text{commitIndex}$ ] then // commit to last log entry
  | if  $t < t^*$  then
  | | return ( $t^*, \text{false}$ )
  | else if log does not match requestor's log at  $\text{prevLogIndex}$  then
  | | return ( $t^*, \text{false}$ )
  | else
  | |  $\text{commitIndex}^* \leftarrow \min(\text{leaderCommit}, \text{index of last new entry})$ 
if  $\text{commitIndex}^* > \text{lastApplied}$  then
  | increment lastApplied, apply  $\text{log}[\text{lastApplied}]$  to state machine
if Timer  $T$  runs out then
  | initialize  $T$  randomly, increment term and start algorithm 7

```

Algorithm 9: Raft Consensus Across-State Responses

362 ► **Lemma 6.** *Algorithm 7 is a correct consensus protocol.*

363 **Proof.** We note that under the definition of $D\&S(V^*)$ a processor will decide upon a value
 364 as soon as it increases its `commitIndex`. Furthermore it will decide upon the value in its first
 365 log entry.

366 We also note that there will never be a majority of processors which failed indefinitely (this
 367 assumption has been made in the original raft algorithm as well). We use this assumption in
 368 order to guarantee that processors will not remain indefinitely in the *wait* commands in the
 369 algorithm. We now prove that the consensus constraints hold.

370 ■ **Validity:** Follows from implementation since the only values proposed for consensus are
 371 taken from the processor's values.

372 ■ **Agreement:** Assume some processor i is the first processor to decide on some value
 373 denoted by v_i . Since the leader is always the first to decide upon a value in a given term,
 374 we may assume i is a leader and denote its term by t_i . Now, assume some other processor
 375 j decided on value v_j during term $t_j \geq t_i$. By the leader completeness property, the entry
 376 $D\&S(v_i)$ will appear in t_j 's leader's log also in the first entry. Thus for j to increase its
 377 `commitIndex` (and decide on a value) it must have accepted t_j 's leader's `AppendEntry`
 378 meaning that by the log matching property their logs must match on at least the first
 379 entry. Therefore j 's first entry would be $D\&S(v_i)$ resulting in $v_j = v_i$.

380 ■ **Termination:** Leader completeness insures us that if someone commits to a value then
 381 eventually all other processors will have that value in their first log index. By the
 382 timing property we are insured that eventually all processors increase their `commitIndex`.
 383 This in turn insures us that eventually all processors will decide upon the same value.
 384 Furthermore, by our assumption that a majority of live processors will eventually exist
 385 ◀

386 We next turn to show how the consensus protocol can be naturally decomposed using
 387 our template. The decomposition works as follows; each term (as described in the consensus
 388 protocol) will now refer to a round in our template. This results in the fact that the protocol
 389 is unending, however eventually all processors will have committed to some value (as in the
 390 original raft protocol).

391 As in our defined VAC the consensus protocol also results in three types of processors.
 392 The first being processors that did not receive a message that a leader was chosen. This
 393 matches the vacillate value of VAC in the sense that they have no guarantee regarding the
 394 state of the system.

395 The second type are processors that received an *AppendEntries* message of the first kind,
 396 i.e. one that **does not** include a change in the commit index. This matches the adopt value
 397 of VAC in that these processors have the guarantee that all other processors which received
 398 such a message received it with the same value (this is ensured by the fact that a majority of
 399 acks is needed in order to send the message).

400 The third and final type are processors that received an *AppendEntries* message of the
 401 second kind, i.e. one that **does** include a change in the commit index. This matches the
 402 commit value of VAC in that these processors are guaranteed that a consensus has been
 403 reached (even if not all processors are aware of it), i.e., that all processors receive the same
 404 value (being accompanied by either an adopt or commit value). This property is ensured by
 405 the *leader completeness* and *state machine safety properties*.

406 The reconciliator in our case is aimed at capturing the timer mechanism. In the consensus
 407 algorithm the timer mechanism was introduced in order to ensure no stalemate was reached,
 408 i.e., to eventually cause convergence. The reconciliator object was introduced to do just that

409 and therefore we define it to capture the timer mechanism as closely as possible. Interestingly
 410 enough in this case, as opposed to Ben-Or for example, it is not the returned value that causes
 411 the wanted behaviour (i.e., prevention of a stalemate) but rather the timing of processors
 412 entering the reconciliator.

413 In algorithm 10 we define our VAC protocol and in algorithm 11 we define the reconciliator
 414 object. We note in addition to the VAC stated we define 2 more changes to the raft consensus
 415 protocol. The first is that if a follower receives an *AppendEntry* message of the first type
 416 (i.e., with an appended entry but without an increase in the commit index) and accepts the
 417 message, then it also sets its X and v values to *adopt* and the value it sees in its last log
 418 entry. The second is that if a follower receives an *AppendEntry* message of the second type
 419 and accepts the message, then it sets X to *commit* and v to the value it sees in its last log
 420 entry.

```

VAC( $v$ )
  ( $X, v^*$ )  $\leftarrow$  (Vacillate,
    log[ $lastLogIndex^*$ ].value)
  state  $\leftarrow$  candidate
  Broadcast RequestVote[ $t^*, id^*$ ,
     $lastLogIndex^*$ ,  $lastLogTerm^*$ ]
  Wait to receive  $> n/2$ 
  ack_RequestVote( $t = t^*$ ,  $granted =$ 
     $true$ )
  Freeze timer T
  ( $X, v^*$ )  $\leftarrow$  (Adopt,
    log[ $lastLogIndex^*$ ].value)
  state  $\leftarrow$  leader
  Broadcast AppendEntries[ $t^*, id^*$ ,
     $prevLogIndex^*$ ,  $prevLogTerm^*$ ,
     $D\&S(v^*)$ ,  $commitIndex^*$ ]
  Wait to receive  $> n/2$ 
  ack_AppendEntries( $t =$ 
     $t^*$ ,  $success = true$ ) /*commitIndex is
    therefore increased, see leader
    responses*/
  ( $X, v^*$ )  $\leftarrow$  (Commit,
    log[ $lastLogIndex^*$ ].value)
  Broadcast AppendEntries[ $t^*, id^*$ ,
     $prevLogIndex^*$ ,  $prevLogTerm^*$ ,
     $-$ ,  $commitIndex^*$ ]
  
```

Algorithm 10: VAC Protocol

```

Reconciliator( $v$ )
  Reset timer and update term
   $D\&S(v) \leftarrow$  log[ $lastLogIndex$ ]
  return  $v$ 
  
```

Algorithm 11: Reconciliator Protocol

421 **► Lemma 7.** *Algorithm 10 is a correct VAC protocol.*

422 **Proof.** In order to ensure that our guarantees hold we prove them for processors which
 423 have not failed during the term. Processors which fail during the term adopt the higher
 424 term once waking up anyhow and should therefore be ignored.

425 **Validity:** Since all values written into the logs were written using received values validity
 426 is insured.

- 427 ■ **Convergence:** We refer the reader to the note following the proof.
- 428 ■ **Termination:** By our assumption that no majority of processors will crash-fail clearly
429 the process will eventually terminate.
- 430 ■ **Coherence over adopt and commit:** As discussed above this is guaranteed by the
431 *leader completeness* and *state machine safety properties*.
- 432 ■ **Coherence over vacillate and adopt:** As discussed above this is ensured by the fact
433 that a majority of acks is needed in order to send an adopt message.
- 434

435 We note that under the raft algorithm infrastructure since consensus is achieved by first
436 electing a leader, convergence does not hold as is. This is indeed plausible due to the fact
437 that the algorithm was made for real world log consistences rather than theoretical consensus.
438 For theoretical purposes one may easily convert the algorithm such that it holds convergence
439 by converting the *wait* steps to *broadcast* steps. I.e., decentralize the messages meaning
440 that instead of electing a leader and having him in charge of logging commands, everyone
441 broadcasts the command they want logged and once someone sees a majority it sends out a
442 commit-to-that-command message. This would result in convergence since if all processors
443 agree on the same value in the first place, all steps would be easily passed.

444 Interestingly enough, this change results in an algorithm that highly resembles Ben-Or's.
445 The only difference is in the way it handles stalemates, or in other words, the reconciliators
446 implemented are different.

447 **5 Adopt-Commit is Not Enough**

448 The concept of decomposing consensus into separate objects is by no means original and was
449 formally presented in [5]. Later work by Aspnes [2] described a framework of *adopt-commit*
450 objects that detect agreement, and *conciliators* that ensure agreement with some probability.
451 We argue that this decomposition fails to capture the inner workings of some well-known
452 algorithms. In these algorithms 3 different types of processors exist throughout the process
453 of achieving consensus; the first are processors which have no guarantee regarding the state of
454 the system. The second are processors that are guaranteed that they are part of a subset
455 of processors that achieved consensus (all other processors are of the first type). The last
456 are processors that are guaranteed that the network has achieved consensus (while not all
457 processors are aware of this fact).

458 In order to make our argument more concrete, we demonstrate how Ben-Or's consensus
459 algorithm cannot be described by a sequence of *adopt-commit* alternating with *conciliator*,
460 while it is naturally described as a sequence of repetitive *vacillate-adopt-commit* followed by
461 *reconciliator*.

To demonstrate the problem with formulating Ben-Or's consensus protocol using

$$U = A_{-1}; A_0; C_1; A_1; C_2; A_2; \dots,$$

462 consider each round of Ben-Or's algorithm [3]². Let P be a processor participating in the
463 agreement process. P experiences one of three possible outcomes: (1) not receiving any *ratify*
464 message. (2) receiving up to t *ratify* messages. (3) receiving more than t *ratify* messages.

² We note that our description follows the presentation of Ben-Or's algorithm given in the survey paper of Aspnes [1]

465 These outcomes correspond to vacillate, adopt, and commit, respectively. Option 1 fits
 466 a processor which received vacillate as it has no guarantees about other values received
 467 by other processors. Option 2 corresponds to *adopt* under the *VAC* framework, since by
 468 *coherence*, any processor that received $(adopt, v)$ is guaranteed that every other processor
 469 that received either *vacillate* or *commit*, also received the value v . Option 3 corresponds to
 470 *commit*, since any processor that received $(commit, v)$ is guaranteed that all other processors
 471 received either $(commit, v)$ or $(adopt, v)$.
 472 However, using only *adopt-commit* objects is not enough in order to describe these three
 473 options.

474 It might be tempting to assume that two consecutive *adopt-commit* objects might resolve
 475 this entanglement as we have shown that *VAC* may be implemented using two *AC* objects.
 476 Note that the concatenation of the *AC* objects (as proposed in [2]) is in a way that is different
 477 than our proposed *VAC* implementation since in their case vacillate-receiving processors are
 478 not represented.

479 We argue this is not the case, that is, we claim that the sequence of $U = A_{-1}; A_0^0; A_0^1; C_1; A_1^0; A_1^1; C_2; \dots$
 480 also fails to describe Ben-Or's consensus protocol. In order to describe option (2) the first
 481 *adopt-commit* must return *adopt* while the second returns *commit*. However, the decomposi-
 482 tion framework described in [2] requires that upon receiving *commit* the processor immediately
 483 decides on the value received, whereas it is possible that in Ben-Or's protocol such a state is
 484 reached with value u but a final agreement is achieved with value $u' \neq u$.

485 6 Conclusions

486 Motivated by the desire to provide a natural decomposition of consensus into building blocks
 487 that describe known algorithms, we defined a more subtle object than *adopt-commit*, the
 488 *vacillate-adopt-commit*, which in turn simplifies the role of the *reconciliator* such that in some
 489 cases it is only a procedure that flips a coin and does not require machinery to ensure *validity*.
 490 Using these building blocks we demonstrate how well known consensus algorithms decompose
 491 into a unified template of a repetitive two step process. We hope a better understanding
 492 of the consensus object may allow research of complexity bounds of the newly introduced
 493 building blocks which in turn may be deduced to consensus.

494 ——— References ———

- 495 1 James Aspnes. Randomized protocols for asynchronous consensus. *Distributed Computing*,
 496 16(2–3):165–175, September 2003.
- 497 2 James Aspnes. A modular approach to shared-memory consensus, with applications to the
 498 probabilistic-write model. *Distributed Computing*, 25(2):179–188, May 2012.
- 499 3 Michael Ben-Or. Another advantage of free choice: Completely asynchronous agreement
 500 protocols (extended abstract). In *Proceedings of the Second Annual ACM SIGACT-SIGOPS*
 501 *Symposium on Principles of Distributed Computing, Montreal, Quebec, Canada, August*
 502 *17-19, 1983*, pages 27–30, 1983. URL: <http://doi.acm.org/10.1145/800221.806707>,
 503 doi:10.1145/800221.806707.
- 504 4 Piotr Berman, Juan A Garay, and Kenneth J Perry. Towards optimal distributed consensus.
 505 In *Foundations of Computer Science, 1989., 30th Annual Symposium on*, pages 410–415.
 506 IEEE.
- 507 5 Eli Gafni. Round-by-round fault detectors (extended abstract): unifying synchrony and
 508 asynchrony. In *Proceedings of the seventeenth annual ACM symposium on Principles of*
 509 *distributed computing*, pages 143–152. ACM, 1998.

23:16 Object Oriented Consensus

- 510 **6** Diego Ongaro and John K. Ousterhout. In search of an understandable consensus algorithm.
511 In *2014 USENIX Annual Technical Conference, USENIX ATC '14, Philadelphia, PA, USA,*
512 *June 19-20, 2014.*, pages 305–319, 2014. URL: [https://www.usenix.org/conference/](https://www.usenix.org/conference/atc14/technical-sessions/presentation/ongaro)
513 [atc14/technical-sessions/presentation/ongaro](https://www.usenix.org/conference/atc14/technical-sessions/presentation/ongaro).
- 514 **7** Marshall Pease, Robert Shostak, and Leslie Lamport. Reaching agreement in the presence
515 of faults. *Journal of the ACM (JACM)*, 27(2):228–234, 1980.