



HAL
open science

Condense and Distill: fast distillation of large floating-point sums via condensation

Stef Graillat, Théo Mary

► **To cite this version:**

Stef Graillat, Théo Mary. Condense and Distill: fast distillation of large floating-point sums via condensation. SIAM Journal on Scientific Computing, In press. hal-04507609v2

HAL Id: hal-04507609

<https://hal.science/hal-04507609v2>

Submitted on 21 Jan 2025

HAL is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

- 38 • The first class relies on the finite number of exponents in typical (IEEE) floating-point
39 arithmetic and on the fact that numbers with exponents not too far apart can be added
40 exactly using extended precision accumulators. This is for example the case of Kulisch’s
41 accumulator [3], which sums all numbers in a very long accumulator, or of the Demmel–Hida
42 algorithm [1], which sums together numbers of comparable exponent using higher precision
43 arithmetic, such as quadruple precision. This strategy is also used in HybridSum [11] and
44 OnlineExactSum [12].
- 45 • The second class exploits the fact that the error incurred by floating-point addition is itself
46 a floating-point number. This error can be computed exactly via error-free transformations
47 such as Fast2Sum. This class includes in particular the AccSum method [8], which computes
48 a faithful rounding of the sum irrespective of its conditioning, and the PrecSum method [9],
49 which computes the sum as if using K -fold precision (for a given K). Those algorithms
50 have been improved respectively as FastAccSum and FastPrecSum in [7].

51 The methods from the first class have a double weakness: they require access to the exponent of
52 the floating-point numbers, which can be expensive, and they require the use of extended precision
53 accumulators. In contrast, the methods from the second class only require standard arithmetic
54 operations on the summands. However, they also require a much larger number of floating-point
55 operations, and their cost strongly depends on the conditioning. Moreover, they are much less
56 parallel than the methods from the first class. In any case, both classes of methods can be quite
57 expensive.

58 Many of the summation methods able to handle ill-conditioned sums rely on the process of
59 *distillation*. Distillation consists in iteratively transforming the original, ill-conditioned sum into an
60 equivalent but well-conditioned sum that can then be evaluated accurately. This is especially the
61 case of the second class of methods mentioned above, although the first class can also be used to
62 design distillation methods, see for example [1, sec. 5] for the Demmel–Hida method.

63 In this article, we propose a method to transform the original sum into another, equivalent
64 sum, which is not necessarily better conditioned, but which is far smaller. The smaller sum can
65 then be distilled inexpensively by any of the distillation methods mentioned above. We call the first
66 transformation of the sum into a smaller equivalent one the process of *condensation*. In the natural
67 language, condensation carries indeed the idea of compacting or contracting a large, complex system
68 into a smaller, simpler one. Moreover, and quite appropriately, in the original physical meaning
69 of the words, condensation (the process of transforming gas to liquid) is a crucial component of
70 the process of distillation (the process of separating substances from a liquid). The entire process
71 (condensing the original sum into a smaller one and then distilling it) thus gives rise to a new
72 method to handle large ill-conditioned sums, which we call the Condense & Distill method.

73 The condensation step is based on the key observation that floating-point numbers with ex-
74 ponents not too far apart can be added exactly, even *in the working precision*, as long as some
75 congruence condition of their least significant digits is satisfied. In particular, base-two numbers
76 with the same exponent and the same least significant bit can be added exactly in the working
77 precision. Condense & Distill thus belongs to the first class of summation methods mentioned
78 above, since it requires access to the exponent of the summands. It is most similar to the Demmel–
79 Hida method: it also adds together numbers of the same exponent, but does not require any
80 extra precision. The entire condensation process can be performed in the working precision; only
81 the final condensed sum needs to be distilled using extended precision (or, for that matter, any
82 other distillation method). This is achieved at the cost of accessing the least significant bit of
83 the summands, which is a negligible overhead compared with the cost of accessing their exponent.

84 Therefore, Condense & Distill allows for significant improvements, not only in terms of performance
85 (because operations in the working precision are faster), but also in terms of robustness/portability.
86 For example, our algorithm can easily accomodate quadruple precision as the working precision.
87 Compared with the second class of summation methods (AccSum, etc.), Condense & Distill shares
88 the main strengths of the first class: its performance can be made completely independent of the
89 conditioning of the sum, and it exhibits nearly perfect parallel scaling.

90 The rest of this article is organized as follows. In [section 2](#), we carry out an analysis to determine
91 conditions for the floating-point addition $x + y$ to be exact. We leverage this analysis in [section 3](#) to
92 develop the Condense & Distill method. We experimentally showcase the use of Condense & Distill
93 against traditional distillation algorithms in [section 4](#), where we also analyze the parallel scaling of
94 our algorithm. Finally we provide our concluding remarks in [section 5](#).

95 **2. When is $x + y$ exact?.** Consider a floating-point number system \mathbb{F} with base β , exponent
96 range (e_{\min}, e_{\max}) , and significand of length $t \geq 2$. To denote $x \in \mathbb{F}$ we will use the notation

$$97 \quad x = \pm(\beta^{e_x} + k_x \varepsilon_{e_x}), \quad \varepsilon_{e_x} = \beta^{e_x+1-t}, \quad k_x \in \mathbb{N},$$

98 where e_x is the unbiased exponent of x and ε_{e_x} is the space between adjacent floating-point numbers
99 in the interval $[\beta^{e_x}, \beta^{e_x+1}]$. Note that $k_x < (\beta - 1)\beta^{t-1}$ since $x < \beta^{e_x+1}$.

100 Given $x, y \in \mathbb{F}$ of the same sign, conditions for the subtraction $x - y$ to be exact are known by
101 Sterbenz lemma (see [\[10\]](#) or [\[6\]](#)). In this section, we determine conditions for the addition $x + y$ to
102 be exact. We begin with the very general result below.

103 **THEOREM 2.1.** *Let $x, y \in \mathbb{F}$ of the same sign $\sigma = \pm 1$ such that*

$$104 \quad x = \sigma(\beta^{e_x} + k_x \varepsilon_{e_x}),$$

$$105 \quad y = \sigma(\beta^{e_y} + k_y \varepsilon_{e_y}).$$

107 *Assuming (without loss of generality) that $|x| \leq |y|$, then $x + y \in \mathbb{F}$, and thus the addition is exact,*
108 *iff one of the following conditions is met:*

- 109 (i) $x = 0$;
110 (ii) $|x + y| < \beta^{e_y+1}$, $e_y - e_x \leq t - 1$, and $k_x \equiv 0 \pmod{\beta^{e_y - e_x}}$;
111 (iii) $|x + y| = \beta^{e_y+1}$, $e_y + 1 \leq e_{\max}$, $e_y - e_x \leq t - 1$, and $k_x \equiv 0 \pmod{\beta^{e_y - e_x}}$;
112 (iv) $|x + y| > \beta^{e_y+1}$, $e_y + 1 \leq e_{\max}$, $e_y - e_x \leq t - 2$, and $k_x + k_y \beta^{e_y - e_x} \equiv 0 \pmod{\beta^{e_y - e_x + 1}}$.

113 *Proof.* Case (i) is trivial. For the remaining cases, we consider positive x and y , the negative
114 case being analagous. We first note that $x + y \in [\beta^{e_y}, \beta^{e_y+2}]$, so that if $x + y$ is to be exact its
115 exponent can only be e_y or $e_y + 1$.

- 116 • Let us first consider the case (ii), where $x + y < \beta^{e_y+1}$. Then $x + y \in \mathbb{F}$ iff ε_{e_y} divides
117 $x + y - \beta^{e_y} = k_y \varepsilon_{e_y} + x$, that is, iff ε_{e_y} divides $x = \beta^{e_x} + k_x \varepsilon_{e_x}$. We first prove that it is
118 necessary for $\varepsilon_{e_y} = \beta^{e_y+1-t}$ to divide β^{e_x} , which is only possible if $e_y - e_x \leq t - 1$. If this
119 condition is not met, then $e_x \leq e_y - t$ and so $\beta^{e_x} \leq \varepsilon_{e_y} / \beta$; moreover, $k_x \varepsilon_{e_x} < (\beta - 1)\beta^{e_x} \leq$
120 $(\beta - 1)\beta^{e_y-t} = (\beta - 1)\varepsilon_{e_y} / \beta$ and therefore $\varepsilon_{e_y} > x$ cannot divide x . We conclude that
121 $e_y - e_x \leq t - 1$ is a necessary condition for $x + y \in \mathbb{F}$ (when $x \neq 0$). The condition
122 is not sufficient, since ε_{e_y} must also divide $k_x \varepsilon_{e_x}$, which happens iff k_x is congruent to
123 0 mod $\beta^{e_y - e_x}$. This concludes case (ii).
- 124 • Case (iii) is identical to case (ii), except we must guarantee that β^{e_y+1} exists by barring
125 overflow with the condition $e_y + 1 \leq e_{\max}$.

126 • In case (iv), $\beta^{e_y+1} < x + y < \beta^{e_y+2}$, we also need $e_y + 1 \leq e_{\max}$ to bar overflow. Then,
127 $x + y \in \mathbb{F}$ iff ε_{e_y+1} divides $x + y - \beta^{e_y+1} = \beta^{e_y}(1 - \beta) + k_y\varepsilon_{e_y} + x$. First, we note that
128 $\varepsilon_{e_y+1} = \beta^{e_y+2-t}$ divides β^{e_y} for $t \geq 2$, so $x + y \in \mathbb{F}$ iff ε_{e_y+1} divides $k_y\varepsilon_{e_y} + \beta^{e_x} + k_x\varepsilon_{e_x}$.
129 Second, we prove that it is necessary for ε_{e_y+1} to divide β^{e_x} , which requires $e_y - e_x \leq t - 2$.
130 Indeed, if this condition is not met, $e_x \leq e_y + 1 - t$, and so $x < \beta^{e_x+1} \leq \beta^{e_y+2-t} = \varepsilon_{e_y+1}$.
131 Thus $x + y < y + \varepsilon_{e_y+1} < \beta^{e_y+1} + \varepsilon_{e_y+1}$, but we are in the case $x + y > \beta^{e_y+1}$, so $x + y \notin \mathbb{F}$.
132 Therefore $e_y - e_x \leq t - 2$ is necessary and ε_{e_y+1} divides β^{e_x} . Moreover, for the condition to
133 be sufficient, we also need ε_{e_y+1} to divide $k_y\varepsilon_{e_y} + k_x\varepsilon_{e_x}$, which happens when $k_x + k_y\beta^{e_y-e_x}$
134 is congruent to $0 \pmod{\beta^{e_y-e_x+1}}$. \square

135 **Theorem 2.1** fully characterizes the conditions for the addition of two floating-point numbers
136 $x + y$ to be exact. The conditions essentially boil down to two components: the exponents of x and
137 y must not be too far apart, and their mantissas must satisfy some congruence condition. Making
138 use of this characterization in practice could however be complex. Interestingly, the conditions
139 become much simpler if we specialize them to numbers sharing the same exponent ($e_x = e_y$).

140 **COROLLARY 2.2.** *Let $x, y \in \mathbb{F}$ of same sign $\sigma = \pm 1$ and same exponent e , such that*

$$141 \quad x = \sigma(\beta^e + k_x\varepsilon_e),$$

$$142 \quad y = \sigma(\beta^e + k_y\varepsilon_e).$$

144 *Then $x + y \in \mathbb{F}$, and thus the addition is exact, iff either*

- 145 (i) $|x + y| < \beta^{e+1}$ or
146 (ii) $e + 1 \leq e_{\max}$ and $k_x + k_y \equiv 0 \pmod{\beta}$.

147 Case (i) of **Corollary 2.2** corresponds to cases (i) and (ii) of **Theorem 2.1**, while case (ii) of the
148 corollary corresponds to cases (iii) and (iv) of the theorem. Note that case (i) of the corollary is only
149 possible for $\beta > 2$, since in binary floating-point arithmetic, the sum of two numbers in the range
150 $(2^e, 2^{e+1}]$ necessarily yields a number in the range $(2^{e+1}, 2^{e+2}]$. Moreover, recall that floating-point
151 numbers can be expressed as

$$152 \quad x = \beta^{e-t} \sum_{i=1}^t d_i \beta^{t-i}, \quad (2.1)$$

153 where the digits d_i satisfy $0 \leq d_i \leq \beta - 1$ with $d_1 \neq 0$ for normalized numbers. Since β divides β^{t-i}
154 for $i < t$, the condition $k_x + k_y \equiv 0 \pmod{\beta}$ further simplifies to a condition on the least significant
155 digits d_t^x and d_t^y of x and y :

$$156 \quad d_t^x + d_t^y \equiv 0 \pmod{\beta}, \quad (2.2)$$

157 that is, $d_t^x + d_t^y$ must be either 0 or β . For $\beta = 2$, this further simplifies to $d_t^x = d_t^y$, yielding the
158 following result.

159 **COROLLARY 2.3.** *If $x, y \in \mathbb{F}$ with $\beta = 2$ have the same sign, exponent, and least significant bit,
160 then barring overflow their addition is exact.*

161 **Corollary 2.3** provides a necessary condition that is much simpler to check, since it only requires
162 to access the sign, exponent, and least significant bit of x and y . In the next section we propose
163 an algorithm that exploits this observation to condense a large sum into an equivalent one with far
164 fewer summands.

165 **3. Fast distillation via condensation.** We now describe the Condense & Distill algorithm,
 166 which exploits [Corollary 2.3](#) to compute rapidly and exactly

167
$$s = \sum_{i=1}^n x_i, \quad x_i \in \mathbb{F}. \quad (3.1)$$

168 Condense & Distill consists of two steps. The first step is to condense the sum by adding pairs
 169 of summands sharing the same exponent, sign, and least significant bit (hereinafter abbreviated as
 170 LSB), until no such pairs remain. As we will prove below, the number of remaining summands is
 171 then bounded by a small value. This first condensation step therefore transforms the original sum
 172 into another sum with a much smaller number of summands. The second step is to then distill this
 173 much smaller sum via any traditional distillation method. The condensation step thus serves as a
 174 preprocessing to accelerate the distillation step.

175 To prove the algorithm’s exactness and cost, we conceptually describe it as building a forest
 176 (a disjoint set of trees). We first place the summands x_i as leaf nodes on a level determined by
 177 their exponent. Then we repeatedly sum pairs of siblings with identical sign and LSB and place the
 178 (exact) result on the level above, until all pairs of siblings have either a different sign or a different
 179 LSB. At this point there thus remains at most 4 nodes per level, and the number of non-empty
 180 levels is itself bounded by $L = \lceil \log_2 n \rceil + d$, where d is a constant independent of n that equals
 181 the number of different exponents among the summands x_i . L is certainly bounded by the total
 182 number of possible exponent values of the floating-point system (e.g., 2047 with IEEE binary64),
 183 and can be much smaller for typical datasets which do not cover the entire exponent range. The
 184 final result is therefore given as the exact unevaluated sum of at most $4L$ floating-point numbers.

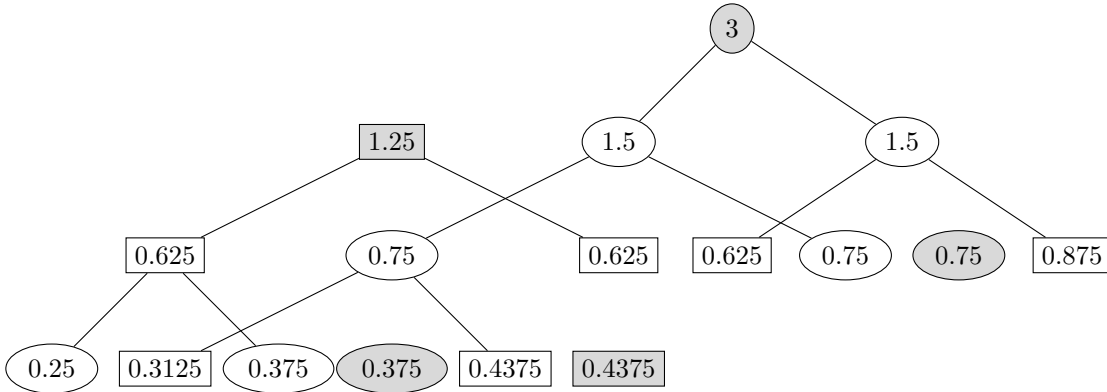


FIG. 3.1. Illustration of the proposed summation algorithm for a simple floating-point system with $t = 3$. The shaded nodes are the remaining values whose unevaluated sum is equal to the exact result. Ellipse and rectangle nodes correspond to numbers with an LSB of 0 and 1, respectively.

185 We illustrate this algorithm in [Figure 3.1](#), using a simple floating-point system \mathbb{F} with $t = 3$.
 186 The eleven leaf nodes correspond to the input summands x_i , whose exact sum is $s = 5.8125$ (we
 187 only consider positive summands here for simplicity, negative numbers would be treated separately
 188 and similarly). All non-leaf nodes correspond to partial results obtained during the computation
 189 and, to easily check that they are indeed floating-point numbers, we provide the list of elements of

190 \mathbb{F} in the interval $(0.25, 3)$:

191 0.25, 0.3125, 0.375, 0.4375, 0.5, 0.625, 0.75, 0.875, 1, 1.25, 1.5, 1.75, 2, 2.5, 3.

192 Ellipse nodes correspond to numbers with an LSB equal to 0, and can thus be summed exactly with
193 other ellipse nodes on the same level. Rectangle nodes correspond to numbers with an LSB equal
194 to 1, and can similarly be summed exactly with other rectangle nodes on the same level. The five
195 shaded nodes are the root nodes and correspond to the remaining numbers that cannot be summed
196 exactly with another node on the same level. They form an unevaluated sum whose result is equal
197 to the exact sum:

198
$$s = 0.375 + 0.4375 + 0.75 + 1.25 + 3.$$

199 It is interesting to remark that the addition of some of these numbers can be represented exactly,
200 namely $0.75 + 1.25 = 2 \in \mathbb{F}$. This is indeed consistent with [Theorem 2.1](#): defining $x = 0.75$ and
201 $y = 1.25$, we have $e_x = -1$, $e_y = 0$, $k_x = 2$, $k_y = 1$, and the condition (iv) of the theorem holds:
202 $k_x + k_y \beta^{e_y - e_x} = 4 \equiv 0 \pmod{2^2}$. However, our algorithm will not exploit this because it only tries
203 to sum numbers of identical exponent, for which we only need to check the LSB.

204 It is important to note that the forest structure of the algorithm is purely conceptual and does
205 not actually need to be built. We also do not need access to all summands previous to the beginning
206 of the computation. [Algorithm 3.1](#) describes an online implementation (which adds summands as
207 they become available, and in any order) that requires at most $4L$ accumulators.

208 We note that the assumption in [Corollary 2.3](#) that x and y have the same sign is not required:
209 if their signs are different, the subtraction $x + y$ is exact by Sterbenz’s lemma, since numbers with
210 the same exponent certainly satisfy $x/2 \leq y \leq 2x$, and this holds regardless of the LSB of x and y .
211 Therefore, the algorithm could also add pairs of summands with the same exponent and different
212 signs, reducing the maximum number of accumulators (and terms in the unevaluated result) from
213 $4L$ to $2L$. However, the drawback is that we would need to recompute the exponent of the result
214 of each addition, since the exponent of a subtraction of two numbers with the same exponent can
215 have an arbitrarily small exponent depending on how close the two numbers are. In contrast, by
216 restricting the pairs to have the same sign, we know that the exponent of $x + y$ is exactly one more
217 than that of x and y .

218 **4. Numerical experiments.** We present a set of numerical experiments to assess the perfor-
219 mance of the Condense & Distill method and its behavior with respect to various parameters such
220 as the dimension n and the condition number κ . We also present a parallel implementation of the
221 method and study its scalability.

222 **4.1. Experimental protocol.** All the experiments were performed on one node of the Olympe
223 supercomputer, equipped with two 18-core Intel Skylake processors. All code was compiled with
224 gfortran version 9.3.0 and with the -O3 optimization flag.

225 We test the methods on ill-conditioned sums randomly generated as follows. Assuming $n =$
226 $2k + 1$ is odd (if n is even we simply add one extra zero summand), we generate k random summands
227 x_i in the range $[10^{-e}, 10^e]$, where e is a fixed parameter that determines the width of the dynamic
228 range of the x_i values; we have used $e = 32$ throughout all experiments. We set another k summands
229 to $-x_i$ and set the last summand to $10^e/\kappa$. Finally we randomly shuffle all summands. The exact
230 sum is equal to $10^e/\kappa$, and its conditioning is of the order of κ .

231 **4.2. Comparison with Demmel–Hida and AccSum.** We begin by comparing the perfor-
232 mance of our new Condense & Distill method with that of the Demmel–Hida and AccSum methods.

Algorithm 3.1 Condense & Distill method.

```
1: Input:  $n$  summands  $x_i$  and a distillation method distill
2: Output:  $s = \sum_{i=1}^n x_i$ 

3: Initialize  $\text{Acc}(e, s, b)$  to 0 for  $e = e_{\min} : e_{\max}$ ,  $s \in \{-1, 1\}$ ,  $b \in \{0, 1\}$ .
4: for all  $x_i$  in any order do
5:    $e = \text{exponent}(x_i)$ 
6:    $s = \text{sign}(x_i)$ 
7:    $b = \text{LSB}(x_i)$ 
8:   insert( $\text{Acc}, x_i, e, s, b$ )
9: end for
10:  $x_{\text{condensed}} = \text{gather}(\text{Acc})$ 
11:  $s = \text{distill}(x_{\text{condensed}})$ 

12: function insert( $\text{Acc}, x, e, s, b$ )
13:   if  $\text{Acc}(e, s, b) = 0$  then
14:      $\text{Acc}(e, s, b) = x$ 
15:   else
16:      $x' = \text{Acc}(e, s, b) + x$ 
17:      $\text{Acc}(e, s, b) = 0$ 
18:      $b' = \text{LSB}(x')$ 
19:     insert( $\text{Acc}, x', e + 1, s, b'$ )
20:   end if
21: end function

22: function  $x_{\text{condensed}} = \text{gather}(\text{Acc})$ 
23:    $i = 0$ 
24:   for all nonzero  $\text{Acc}(e, s, b)$  do
25:      $i = i + 1$ 
26:      $x_{\text{condensed}}(i) = \text{Acc}(e, s, b)$ 
27:   end for
28: end function
```

233 [Figure 4.1](#) plots the time cost of each method for varying κ and for two fixed values of n , $n = 10^7$
234 (left) or $n = 10^8$ (right). The figure shows that, as expected, the cost of the Condense & Distill
235 and Demmel–Hida methods is independent of κ , with Condense & Distill being roughly 35% faster
236 because it avoids using quadruple precision. In contrast, the cost of AccSum strongly depends
237 on κ , growing at a rate of roughly $\log \kappa$. As a result, the time comparison between AccSum and
238 Condense & Distill depends on κ : for only moderately ill-conditioned sums, AccSum is faster, but
239 as κ increases Condense & Distill (and even Demmel–Hida) eventually outperforms AccSum, po-
240 tentially by very large factors if the sum is extremely ill conditioned. The cutoff value of κ for
241 which Condense & Distill outperforms AccSum also seems to decrease as n increases: it is equal to
242 $\kappa \simeq 10^{35}$ for $n = 10^7$ and $\kappa \simeq 10^{20}$ for $n = 10^8$.

243 We confirm this last trend in [Figure 4.2](#), where we compare the performance of the methods for
244 an increasing n and a fixed value of κ . The figure shows that Condense & Distill becomes more and

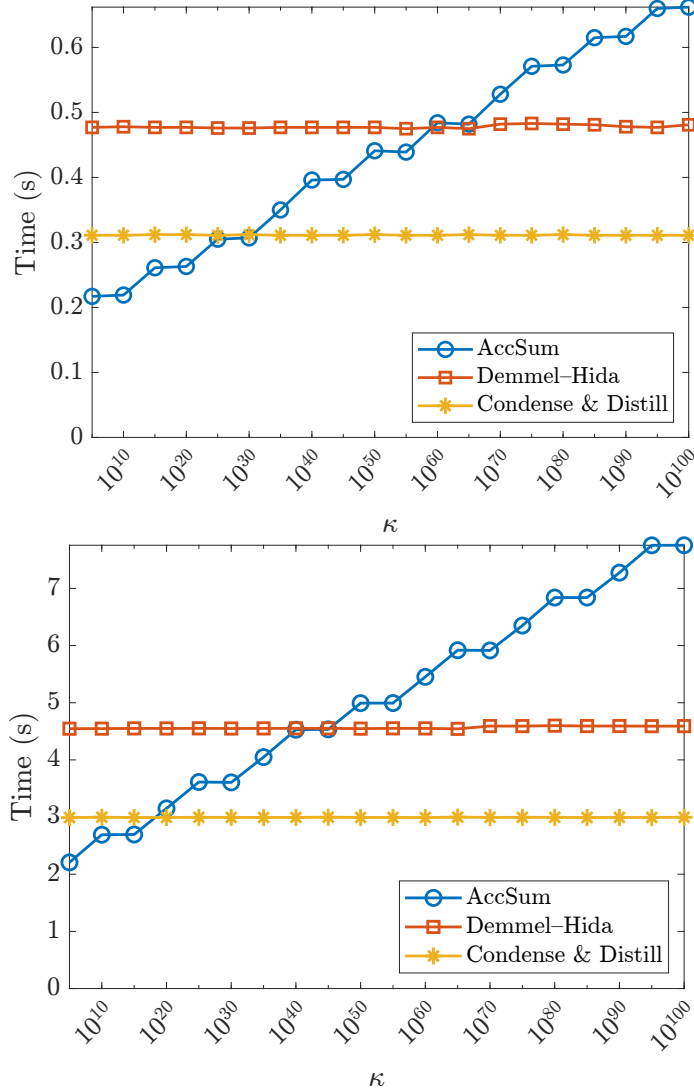


FIG. 4.1. Comparison between the Demmel-Hida, AccSum, and Condense & Distill algorithms, as a function of the condition number κ and for two dimensions $n = 10^7$ (top) and $n = 10^8$ (bottom). All algorithms are run sequentially (1 thread).

245 more competitive with respect to AccSum as n increases. While AccSum is faster for small sums, it
 246 is eventually outperformed by Condense & Distill and even by Demmel-Hida, for sufficiently large
 247 sums. The cutoff value of n for which Condense & Distill outperforms AccSum similarly decreases
 248 as κ increases: for example, it is equal to $n \simeq 10^7$ for $\kappa = 10^{30}$ and $n \simeq 10^6$ for $\kappa = 10^{60}$.

249 **4.3. Parallel scaling.** In the previous comparison, the methods are executed sequentially
 250 (using only 1 thread) but, as mentioned, Condense & Distill, like Demmel-Hida and similar meth-

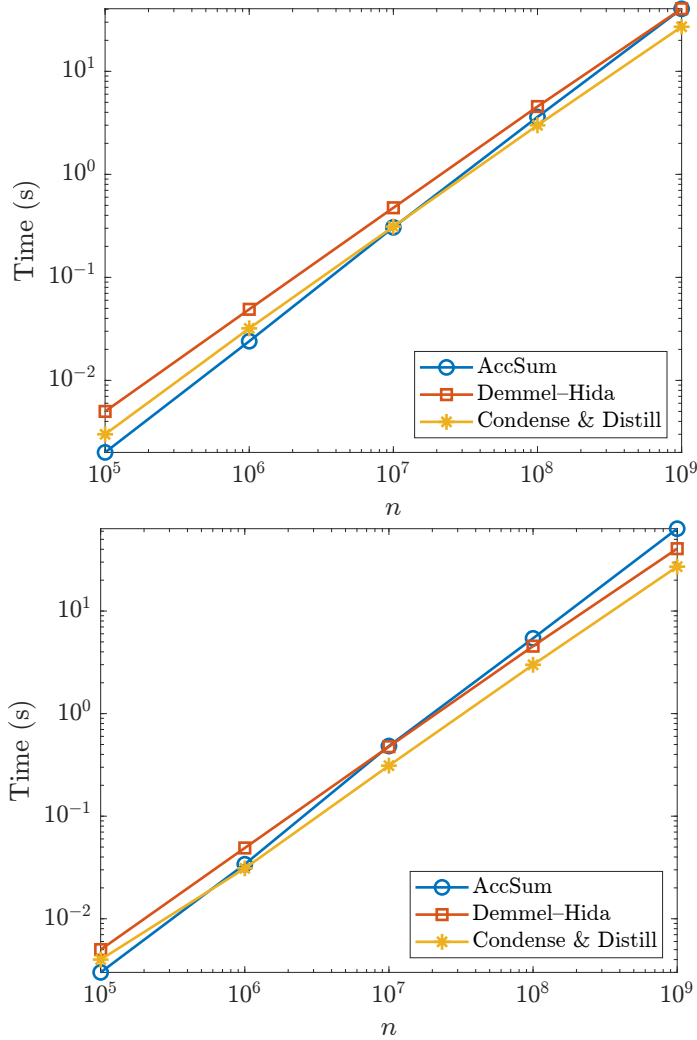


FIG. 4.2. Comparison between the Demmel–Hida, AccSum, and Condense & Distill algorithms, as a function of the number of summands n and for two condition numbers $\kappa = 10^{30}$ (top) and $\kappa = 10^{60}$ (bottom). All algorithms are run sequentially (1 thread).

251 ods, is very amenable to parallelism. We have implemented a parallel version of Condense & Distill,
 252 which can exploit p threads by splitting the summands into p blocks and condensing each block
 253 in parallel. This yields a condensed sum with at most $4Lp$ summands, which can be sequentially
 254 condensed into an even smaller sum with at most $4L$ summands, before being sequentially distilled.
 255 Figure 4.3 analyzes the parallel scaling of this method with a varying number of threads from 1 to
 256 36. We consider both strong scaling (right plot, with fixed $n = 10^8$) and weak scaling (left plot,
 257 with $n = 3 \times 10^6$ per thread). The method exhibits near perfect scaling, as expected.

258 As mentioned, similar scaling can also be expected from the Demmel–Hida method. In contrast,

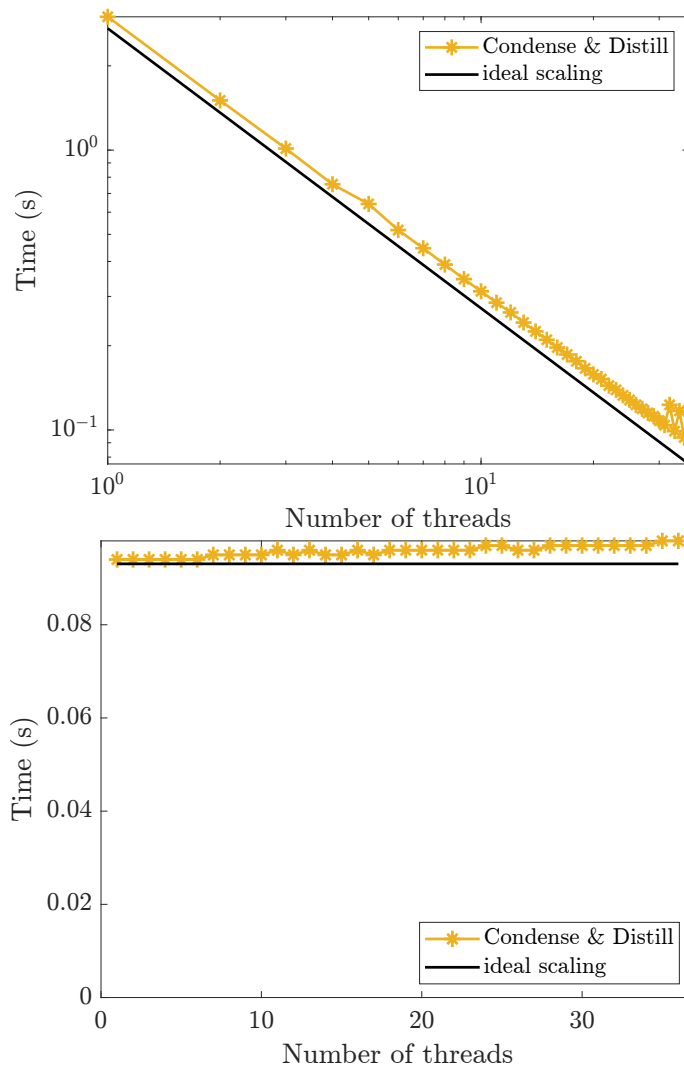


FIG. 4.3. Parallel scaling of Condense & Distill, using from 1 to 36 threads. Top: strong scaling ($n = 10^8$ is fixed). Bottom: weak scaling ($n = 3 \times 10^6$ per thread).

259 AccSum and similar methods offer much less parallelism. We do not study the parallel scaling of
 260 AccSum here, but refer to [5], which shows that AccSumK can achieve at best a parallel efficiency
 261 of only 50%. Therefore, in a parallel setting, we can expect the performance comparison between
 262 Condense & Distill and AccSum to be even more in favor of the former, even for small values of κ .

263 **4.4. Quadruple precision as the working precision.** We finally illustrate how condensa-
 264 tion can be even more beneficial in the case where the working precision is quadruple precision,
 265 which may for example be necessary for applications requiring a high level of accuracy. In this
 266 situation Condense & Distill is clearly at an advantage compared with both Demmel–Hida and

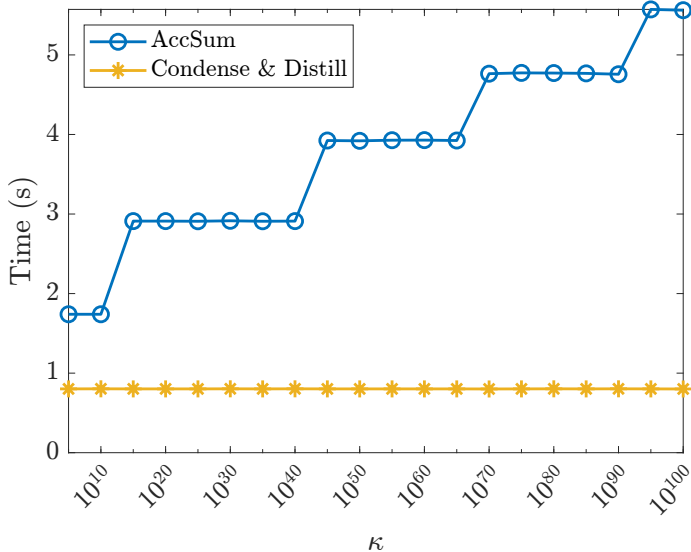


FIG. 4.4. Comparison between AccSum and Condense & Distill using quadruple precision as the working precision ($n = 10^7$, both algorithms are executed using 1 thread).

267 AccSum. Indeed, since Demmel–Hida requires extended precision, using quadruple precision as
 268 the working precision would require access to an even higher precision, which is unavailble on most
 269 architectures. The only solution would be to rely on an arbitrary precision library, but this would
 270 likely be very expensive and we do not explore this option further. As for AccSum, it can easily be
 271 executed in quadruple precision since it also runs entirely in the working precision. However, the
 272 cost comparison with Condense & Distill tips even more in favor of the latter, because the relative
 273 cost of accessing the summands exponent is smaller with respect to the cost of arithmetic operations
 274 (which are much more expensive in quadruple precision). This is illustrated in Figure 4.4, which
 275 shows that Condense & Distill achieves even larger speedups with respect to AccSum, and even for
 276 small condition numbers.

277 **5. Conclusion.** We have proposed a new distillation method, Condense & Distill (Algo-
 278 rithm 3.1), which employs a preprocessing condensation step to turn a large ill-conditioned sum
 279 into a still ill-conditioned, but far smaller sum, which is then distilled inexpensively via traditional
 280 distillation methods. The condensation step relies on Corollary 2.3, which proves that floating-point
 281 numbers with the same exponent and least significant bit can be added exactly. Compared with
 282 other summation methods that also require accessing the exponent field of the summands, such as
 283 the Demmel–Hida method [1], Condense & Distill can run entirely in the working precision. As a
 284 result, Condense & Distill is faster, and is also more portable since it does not require any extra
 285 precision to be available. Compared with distillation methods based on error-free transformations,
 286 such as AccSum [8], Condense & Distill’s cost does not increase with the conditioning, and exhibits
 287 much better parallel scaling. Overall, we have thus shown Condense & Distill to be an efficient
 288 method to distill large ill-conditioned sums.

289 **Acknowledgments.** We thank Massimiliano Fasi and Mantas Mikaitis for a discussion at the
290 2022 Creativity workshop in Manchester, organized by Nick Higham and Dennis Sherwood, that
291 led to the observation that numbers with the same exponent and least significant bit can be added
292 exactly.

293 **Funding.** This work was partially supported by the InterFLOP (ANR-20-CE46-0009), NuS-
294 CAP (ANR-20-CE48-0014), and MixHPC (ANR-23-CE46-0005-01) projects of the French National
295 Agency for Research (ANR).

296

REFERENCES

- 297 [1] J. DEMMEL AND Y. HIDA, *Accurate and efficient floating point summation*, SIAM Journal on Scientific Com-
298 puting, 25 (2004), pp. 1214–1248.
- 299 [2] N. J. HIGHAM, *Accuracy and Stability of Numerical Algorithms*, Society for Industrial and Applied Mathematics,
300 Philadelphia, PA, USA, second ed., 2002.
- 301 [3] U. W. KULISCH AND W. L. MIRANKER, *The arithmetic of the digital computer: A new approach*, SIAM Review,
302 28 (1986), pp. 1–40.
- 303 [4] M. LANGE, *Toward accurate and fast summation*, ACM Trans. Math. Softw., 48 (2022).
- 304 [5] X. LEI, T. GU, S. GRAILLAT, H. JIANG, AND J. QI, *A fast parallel high-precision summation algorithm based*
305 *on accsumk*, Journal of Computational and Applied Mathematics, 406 (2022), p. 113827.
- 306 [6] J.-M. MULLER, N. BRUNIE, F. DE DINECHIN, C.-P. JEANNEROD, M. JOLDES, V. LEFÈVRE, G. MELQUIOND,
307 N. REVOL, AND S. TORRES, *Handbook of floating-point arithmetic*, Birkhäuser/Springer, Cham, second ed.,
308 2018.
- 309 [7] S. M. RUMP, *Ultimately fast accurate summation*, SIAM Journal on Scientific Computing, 31 (2009), pp. 3466–
310 3502.
- 311 [8] S. M. RUMP, T. OGITA, AND S. OISHI, *Accurate floating-point summation part i: Faithful rounding*, SIAM
312 Journal on Scientific Computing, 31 (2008), pp. 189–224.
- 313 [9] S. M. RUMP, T. OGITA, AND S. OISHI, *Fast high precision summation*, Nonlinear Theory and Its Applications,
314 IEICE, 1 (2010), pp. 2–24.
- 315 [10] P. H. STERBENZ, *Floating-point computation*, Prentice-Hall Inc., Englewood Cliffs, N.J., 1974. Prentice-Hall
316 Series in Automatic Computation.
- 317 [11] Y.-K. ZHU AND W. B. HAYES, *Correct rounding and hybrid approach to exact floating-point summation*, SIAM
318 J. Sci. Comput., 31 (2009), pp. 2981–3001.
- 319 [12] Y.-K. ZHU AND W. B. HAYES, *Algorithm 908: Online exact summation of floating-point streams*, ACM Trans-
320 actions on Mathematical Software (TOMS), 37 (2010), pp. 1–13.