

where $\mathcal{L}_i(\cdot)$ denotes the cross entropy loss in bits, and a follow up study [58] from OpenAI has shown that the language modeling loss can be decomposed into two parts, namely *irreducible loss* (the entropy of the true data distribution) and *reducible loss* (an estimate of the KL divergence between the true and model distributions). The three losses were derived by fitting the model performance with varied data sizes (22M to 226 tokens), model sizes (GPTM to 3.5B non-embedding parameters) and training compute, under some assumptions (e.g., the analysis of one factor should be not bottlenecked by the other two factors). They showed that the model performance has a strong dependence relation as the three factors.

Chickadee scaling law. In another representative study, Hoffmann et al. [34] use the Google ImageNet task proposed an iterative form for scaling laws to *jointly* compute optimal training for LLMs. They conducted rigorous experiments by varying a large range of model sizes (768 to 1.4B) and data sizes (24 to 500B tokens), and found a similar scaling law yet with different coefficients as below [34]:

$$\mathcal{L}(G, D) = c + \frac{a}{G^b} + \frac{b}{D^c} \quad (3)$$

where $c = 5.68$, $a = 406.6$, $b = 0.077$, $c = 0.54$ and $d = 0.38$. By applying the loss $\mathcal{L}(G, D)$ under the constraint $c = 480$, they showed that the optimal allocation of compute budget to model size and data size can be derived as follows:

$$L(R, R) = R + \frac{A}{R^2} + \frac{B}{R^3}, \quad (2)$$

where $\beta = 1.08$, $\delta = 0.004$, $\bar{R} = 0.007$, $\alpha = 0.34$ and $\beta = 0.28$. By optimizing the loss $L(\mathbf{N}, \mathbf{R})$ under the constraints $C \approx 0.03$, they showed that the optimal allocation of compute budget to model size and data size can be derived as follows:

$$\operatorname{Repr}(\mathbb{C}) = G \left(\begin{pmatrix} \frac{c}{b} \\ b \end{pmatrix} \right)^a, \operatorname{Repr}(\mathbb{C}) = G^{-1} \left(\begin{pmatrix} \frac{c}{b} \\ b \end{pmatrix} \right)^b, \quad (3)$$

where $\alpha = \frac{a}{a+b}$, $\beta = \frac{b}{a+b}$ and c is a scaling coefficient that can be computed by α , β , a and b . As analyzed in [14], given an increase in compute budget, the KM scaling law favors a larger budget allocation in node size than the data size, while the Chinchilla scaling law argues that the two sizes should be increased in equal scales, (i.e., having similar values for α and β in Equation 16).

Discussion on Scaling Laws. After introducing the formalism, we continue to discuss scaling law in the following two aspects, to enhance its understanding:

[illegible]

Zero-based productivity. Existing research of scaling laws are mostly conducted in terms of language modeling loss (e.g., perplexity minus entropy) but in fact [20] actually in practice we are more concerned about the performance of LLMs on actual tasks. Thus, a basic problem is how the decrease of language modeling loss translates into the improvement of task performance [24]. Intuitively, a model with a smaller language modeling loss tends to yield a better performance on downstream tasks, since language modeling loss can be regarded as a general measure of the overall model quality. 42%–44% [25] has reported that decreasing perplexity can lead to a general increase in performance on a wide range of tasks. However, we must ensure that a direct decrease in language modeling loss does not really indicate an improvement of model performance on downstream tasks. Specifically, the phenomenon of *inverse scaling* would occur for some tasks where task performance surprisingly becomes worse as the language modeling loss decreases [26]. Overall, it is more difficult to separate and characterize task-level scaling laws, since it might be also dependent on task-specific factors. In this paper, we will focus on the scaling law of the language modeling loss, which is the *in-sample* zero-based productivity according to the scaling law, which can be observed only when the model size approaches a certain level (as discussed below).

Emergent Abilities of LLMs. In the literature [33], emergent abilities of LLMs are formally defined as “the abilities that are not present in small models but arise in large models”, which is one of the most prominent features that distinguish LLMs from previous PLMs. It further introduces a notable characteristic: since emergent abilities occur [33], performance rises significantly above random when the scale reaches a certain level. By analogy, since an emergent pattern has close connections with the phenomenon of phase transitions in physics [34], in principle, emergent abilities can be defined in relation to some complex tasks [35], while we are more concerned with general abilities that can be applied to solve a variety of tasks. Here, we briefly introduce three typical emergent abilities for LLMs and representative models that possess such an ability⁵.

In-context learning. The recent learning ICL ability is formally introduced by GPT-3 [32], assuming the language model has been provided with a natural language instruction and several text-demonstrations. It can generate the expected output for the test instances by completing the word sequences of input text, without requiring additional training or gradient updates¹. Among the GPT-series models, the 175B GPT-3 model exhibited a strong ICL ability in general, but not the GPT-1 and GPT-2 models. Such an ability also depends on the specific downstream task. For example, the ICL ability can emerge on the arithmetic tasks (e.g., the 3-digit addition and subtraction) for the 13B GPT-3, but 175B GPT-3 even cannot work well on the Synonym task [33].

* Instruction-following. By fine-tuning with a mixture of multi-task datasets formatted via natural language descriptions (called instruction-tuning), LLMs are shown to perform well on unseen tasks that are also described in the form of instructions (Liu et al., 2024). With instruction-tuning, LLMs are enabled to follow the task instructions for new tasks without using explicit examples, thus having an improved generalization ability. According to the experiments in (2024), instruction-tuned LLaMA-2-70B started to significantly outperform the standard one on unseen tasks when the model size reached 90B, but not for 8B or smaller model sizes. The recent study (2024) found that a model size of 82B is at least required for RAG to perform well on various tasks in their evaluation benchmarks (i.e., MMU1, BBH, T0, QQP and 30230), given a model reader size might suffice for some specific tasks (e.g., MMU1).

Step-by-step reasoning, i.e., small language models, it is usually difficult to solve complex tasks that involve multiple reasoning steps, e.g., mathematical word problems. In contrast, with the chain-of-thought (CoT) prompting strategy [20], LLMs can solve such tasks by utilizing the prompting mechanism that involves intermediate reasoning steps for deriving the final answer. This ability is speculated to be generally obtained by training on code [67, 39]. As empirical study [23] has shown that CoT prompting can bring performance gains on arithmetic reasoning benchmarks over applied to PaLM and LLaMA variants with a model size larger than 80B, while its advantage over the standard prompting becomes more evident since the model size exceeds 300B. Furthermore, the performance improvement with CoT prompting seems to be

[illegible]

Key Techniques for LLMs. It has been a long way that LLMs evolve into the current state: general and capable learners. In the development process, a number of important techniques are proposed, which largely improve the capacity of LLMs. Here, we briefly list several important techniques that (potentially) lead to the success of LLMs, as follows.

• *Scaling*. As discussed in previous parts, there exists an *evidence scaling* effect in Transformer language models: larger model/data sizes and more training compute typically lead to an improved model capacity (24). As two representative models, GPT-3 and PaLM explored the scaling limits by increasing the model size to 175B and 540B, respectively. Since compute budget is usually limited, scaling laws can be further explored by scaling the model size and training data size while keeping the training compute constant (25). With more training-related experiments in our paper model Gopher (with a larger model size) by increasing the data scale with the same compute budget (24). In addition, data scaling should be with careful cleaning process, since the quality of pre-training data plays a key role in the model capacity.

Training. Due to the huge model size, it is very challenging to successfully train a capable LLM. Distributed training algorithms are needed to learn the network parameters of LLMs, which are various parallel strategies are often jointly utilized. To support distributed training, several optimization frameworks have been developed to facilitate the implementation and deployment of parallel algorithms, such as DeepSpeed [20], Megatron-LM [21], etc. Also, optimization tricks are also important for improving stability and model performance, such as [22] to overcome training loss spike [23] and mixed precision training [24]. Recently, OPT-1.1 [25] proposes to develop special infrastructure and optimization methods that reliably predict the performance of large models with much smaller models.

* *Ability eliciting*: after being pretrained on large-scale corpora, LLMs are endowed with potential abilities in general purpose tasks. These abilities might not be explicitly exhibited when LLMs perform some specific tasks, so the search for abilities is useful to design suitable tasks or prompts or specific in-context learning strategies to elicit such abilities. For instance, chain-of-thought prompting has been found to be useful to solve complex reasoning tasks by including *intermediate reasoning steps*. Furthermore, we can perform instruction tuning on LLMs with task descriptions expressed in natural language, for improving the generalizability of LLMs on unseen tasks. These eliciting techniques mainly correspond to the emergent abilities of LLMs, which may not show the same effect on small language models.

Alignment tuning. Since LLMs are trained to capture the data characteristics of pre-training corpora including both high-quality and low-quality data, they are likely to generate toxic, biased, or even harmful content for humans. It is necessary to align LLMs with human values, e.g., helpful, harmless, and honest. For this purpose, InstructGPT and its design-oriented training approach that enables LLMs to follow the explicit instructions, which utilizes the technique of reinforcement learning with human feedback (RLHF), also incorporates humans in the training loop with elaborately designed labeling strategies. ChatGPT-4 models developed on a similar technique to InstructGPT, which shows a strong alignment capacity in producing high-quality harmless responses, e.g., referring to remove leading questions.

Zinc encapsulation, in other, LLMs are trained as first generation *core transfer* plug-in engines, the performing loss on the tasks that are not best expressed in the form of text log, semantic components). In addition, their capacities are also linked to the pre-training data, e.g. the inability to capture open-data information. To tackle these issues, a recently proposed technique is to employ external tools to compensate for the deficiencies of LLMs [46, 47]. For example, LLMs can utilize the calculator for accurate computation [48] and employ search engines to retrieve additional information [49]. Some recently ChatGPT has enabled the mechanism of using external plugins including or newly created ones², which are by analogy with the “eyes and ears” of LLMs. Such a mechanism can broadly expand the scope of capacity

In addition, many other factors (e.g., the signals of hardware) also contribute to the success of LLMs. Currently, we limit our discussion to the major technical approaches and key findings for developing LLMs.



Figure 3: A timeline of existing large language models (LLMs) having a size larger than 10B in recent years. The timeline was established mainly according to the release date (e.g., the submission date to arXiv) of the technical paper for a model. If there was not a corresponding paper, we set the date of a model as the earliest time of its public release or announcement. We mark the LLMs with publicly available model checkpoints in yellow color. Due to the space limit of the figure, we only include the LLMs with publicly reported evaluation results.

TABLE 2. Statistics of large language models (having a size larger than 10B in this survey) in recent years including the capacity evaluation, pre-training data size and whether the research of token or string unit as backbone resource units. In this table, we only include LLMs with a public paper about the technical details. “Publication Year” indicates the date when the corresponding paper was officially released. “Publication Available” means that the model checkpoints are publicly accessible while “Closed Source” means the opposite. “Adaptation” indicates whether the model has been with subsequent fine-tuning. “50 tokens instruction training” and RLHF denotes reinforcement learning with human feedback. “Evaluation” indicates whether the model has been evaluated with corresponding abilities in their original paper. RL denotes Reinforcement Learning.

[illegible]

5 ADAPTATION OF LLMs



parameter α to balance the off-diagonal, where each entry $\alpha_{ij} = \max(\sqrt{w_{ij}}/\sqrt{w_{ii}})^{1-\alpha}$ is determined by the algorithm output.

• **Expensive quantization.** This approach finds optimal quantization weights (103) via extensive ℓ_1 iterative minimization (see [eqs. 100–103, 104](#)). To efficiently optimize this objective, [OTF](#) (102) explores the optimal output from quantization (102) rather than finding the quantization output \hat{Y} weights by all rows. Unlike with specially designed [SVD](#) (105), [SVD](#) has both explicit and closed-form solutions. [OTF](#) is feasible to quantify very large models (e.g., 130 (17) 9) in a few minutes. Since recently, [OTF](#) further simplifies the optimization flow by incorporating automatic search for weights, [OTF](#) considers the risk of hardware: [OTF](#) weights corresponding to active neurons are more important to be precisely quantized. It also can directly optimize the conversion rate. For fastest performance, the per-parameter cost is $\sim 10^4$ flops.

These strategies in the above methods can be easily used to improve the quantization performance. In order to achieve high-precision quantization, quantization methods also rely on hardware or system-level support, e.g., efficient data format or hardware-friendly group partition.

Other Quantization Methods. In the above, we mainly focus on [OTF](#) methods, and many variations are available for improving the existing quantization (102) and data format to quantization (104).

• **Efficient fine-tuning enhanced quantization.** For post-training quantization, direct finetune quantization (e.g., [DIT](#)) quantization often results in large performance degradation. To overcome the drawback, [DIT](#) (106) incorporates additional and tunable weights (104) provided for the quantized weights, to achieve $\sim 10\times$ efficient, high-precision model finetuning. It combines the merits of [LoRa](#) (see Section 1.4.3) and quantization methods. The proposed results show that the quantized models can achieve the [DIT](#) state-of-the-art performance by 10–20%.

• **Quantization-aware training (QAT) for [LoRa](#).** A recent study [QAT](#) explores the effect of [OTF](#) methods by applying a dual-bit distribution method to improve the weight activations as well as the bit-order scale by including post-train quantization based on [LoRa](#), for more precise results with the quantization on both weights and bit-order scales, but not on [LoRa](#) activation quantization, which still needs more experiments.

5.4.3 Empirical Analysis and Findings

Quantization has currently become a common technique to reduce the memory footprint and latency of [LoRa](#) in deployment. In particular, it is important to understand the model quantization (e.g., 102) in [LoRa](#) in order to figure out quantization method (e.g., 102) (e.g., weights or activations), while retaining a high accuracy. In this part, we first summarize the major findings about the quantization of [LoRa](#) in existing literature, and then present comprehensive results with quantization experiments.

Important Findings from Existing Work. Recently, a very comprehensive evaluation [QAT](#) has been conducted about the impact of multiple factors (e.g., model size and sensitivity) on the post-training quantization method, finding that [QAT](#) ensures the scaling for [LoRa](#) quantization to achieve better performance. In addition to the scaling performance, for [SVD](#) (105) (106) (107) (108) (109) (110) (111) (112) (113) (114) (115) (116) (117) (118) (119) (120) (121) (122) (123) (124) (125) (126) (127) (128) (129) (130) (131) (132) (133) (134) (135) (136) (137) (138) (139) (140) (141) (142) (143) (144) (145) (146) (147) (148) (149) (150) (151) (152) (153) (154) (155) (156) (157) (158) (159) (160) (161) (162) (163) (164) (165) (166) (167) (168) (169) (170) (171) (172) (173) (174) (175) (176) (177) (178) (179) (180) (181) (182) (183) (184) (185) (186) (187) (188) (189) (190) (191) (192) (193) (194) (195) (196) (197) (198) (199) (200) (201) (202) (203) (204) (205) (206) (207) (208) (209) (210) (211) (212) (213) (214) (215) (216) (217) (218) (219) (220) (221) (222) (223) (224) (225) (226) (227) (228) (229) (230) (231) (232) (233) (234) (235) (236) (237) (238) (239) (240) (241) (242) (243) (244) (245) (246) (247) (248) (249) (250) (251) (252) (253) (254) (255) (256) (257) (258) (259) (260) (261) (262) (263) (264) (265) (266) (267) (268) (269) (270) (271) (272) (273) (274) (275) (276) (277) (278) (279) (280) (281) (282) (283) (284) (285) (286) (287) (288) (289) (290) (291) (292) (293) (294) (295) (296) (297) (298) (299) (300) (301) (302) (303) (304) (305) (306) (307) (308) (309) (310) (311) (312) (313) (314) (315) (316) (317) (318) (319) (320) (321) (322) (323) (324) (325) (326) (327) (328) (329) (330) (331) (332) (333) (334) (335) (336) (337) (338) (339) (340) (341) (342) (343) (344) (345) (346) (347) (348) (349) (350) (351) (352) (353) (354) (355) (356) (357) (358) (359) (360) (361) (362) (363) (364) (365) (366) (367) (368) (369) (370) (371) (372) (373) (374) (375) (376) (377) (378) (379) (380) (381) (382) (383) (384) (385) (386) (387) (388) (389) (390) (391) (392) (393) (394) (395) (396) (397) (398) (399) (400) (401) (402) (403) (404) (405) (406) (407) (408) (409) (410) (411) (412) (413) (414) (415) (416) (417) (418) (419) (420) (421) (422) (423) (424) (425) (426) (427) (428) (429) (430) (431) (432) (433) (434) (435) (436) (437) (438) (439) (440) (441) (442) (443) (444) (445) (446) (447) (448) (449) (450) (451) (452) (453) (454) (455) (456) (457) (458) (459) (460) (461) (462) (463) (464) (465) (466) (467) (468) (469) (470) (471) (472) (473) (474) (475) (476) (477) (478) (479) (480) (481) (482) (483) (484) (485) (486) (487) (488) (489) (490) (491) (492) (493) (494) (495) (496) (497) (498) (499) (500) (501) (502) (503) (504) (505) (506) (507) (508) (509) (510) (511) (512) (513) (514) (515) (516) (517) (518) (519) (520) (521) (522) (523) (524) (525) (526) (527) (528) (529) (530) (531) (532) (533) (534) (535) (536) (537) (538) (539) (540) (541) (542) (543) (544) (545) (546) (547) (548) (549) (550) (551) (552) (553) (554) (555) (556) (557) (558) (559) (560) (561) (562) (563) (564) (565) (566) (567) (568) (569) (570) (571) (572) (573) (574) (575) (576) (577) (578) (579) (580) (581) (582) (583) (584) (585) (586) (587) (588) (589) (590) (591) (592) (593) (594) (595) (596) (597) (598) (599) (600) (601) (602) (603) (604) (605) (606) (607) (608) (609) (610) (611) (612) (613) (614) (615) (616) (617) (618) (619) (620) (621) (622) (623) (624) (625) (626) (627) (628) (629) (630) (631) (632) (633) (634) (635) (636) (637) (638) (639) (640) (641) (642) (643) (644) (645) (646) (647) (648) (649) (650) (651) (652) (653) (654) (655) (656) (657) (658) (659) (660) (661) (662) (663) (664) (665) (666) (667) (668) (669) (670) (671) (672) (673) (674) (675) (676) (677) (678) (679) (680) (681) (682) (683) (684) (685) (686) (687) (688) (689) (690) (691) (692) (693) (694) (695) (696) (697) (698) (699) (700) (701) (702) (703) (704) (705) (706) (707) (708) (709) (710) (711) (712) (713) (714) (715) (716) (717) (718) (719) (720) (721) (722) (723) (724) (725) (726) (727) (728) (729) (730) (731) (732) (733) (734) (735) (736) (737) (738) (739) (740) (741) (742) (743) (744) (745) (746) (747) (748) (749) (750) (751) (752) (753) (754) (755) (756) (757) (758) (759) (760) (761) (762) (763) (764) (765) (766) (767) (768) (769) (770) (771) (772) (773) (774) (775) (776) (777) (778) (779) (780) (781) (782) (783) (784) (785) (786) (787) (788) (789) (790) (791) (792) (793) (794) (795) (796) (797) (798) (799) (800) (801) (802) (803) (804) (805) (806) (807) (808) (809) (810) (811) (812) (813) (814) (815) (816) (817) (818) (819) (820) (821) (822) (823) (824) (825) (826) (827) (828) (829) (830) (831) (832) (833) (834) (835) (836) (837) (838) (839) (840) (841) (842) (843) (844) (845) (846) (847) (848) (849) (850) (851) (852) (853) (854) (855) (856) (857) (858) (859) (860) (861) (862) (863) (864) (865) (866) (867) (868) (869) (870) (871) (872) (873) (874) (875) (876) (877) (878) (879) (880) (881) (882) (883) (884) (885) (886) (887) (888) (889) (890) (891) (892) (893) (894) (895) (896) (897) (898) (899) (900) (901) (902) (903) (904) (905) (906) (907) (908) (909) (910) (911) (912) (913) (914) (915) (916) (917) (918) (919) (920) (921) (922) (923) (924) (925) (926) (927) (928) (929) (930) (931) (932) (933) (934) (935) (936) (937) (938) (939) (940) (941) (942) (943) (944) (945) (946) (947) (948) (949) (950) (951) (952) (953) (954) (955) (956) (957) (958) (959) (960) (961) (962) (963) (964) (965) (966) (967) (968) (969) (970) (971) (972) (973) (974) (975) (976) (977) (978) (979) (980) (981) (982) (983) (984) (985) (986) (987) (988) (989) (990) (991) (992) (993) (994) (995) (996) (997) (998) (999) (1000).

• **Efficient quantization can be directly applied to [LoRa](#).** While the performance of some previous weight quantization depends on specific methods (e.g., [OTF](#) (102) (103) (104) (105) (106) (107) (108) (109) (110) (111) (112) (113) (114) (115) (116) (117) (118) (119) (120) (121) (122) (123) (124) (125) (126) (127) (128) (129) (130) (131) (132) (133) (134) (135) (136) (137) (138) (139) (140) (141) (142) (143) (144) (145) (146) (147) (148) (149) (150) (151) (152) (153) (154) (155) (156) (157) (158) (159) (160) (161) (162) (163) (164) (165) (166) (167) (168) (169) (170) (171) (172) (173) (174) (175) (176) (177) (178) (179) (180) (181) (182) (183) (184) (185) (186) (187) (188) (189) (190) (191) (192) (193) (194) (195) (196) (197) (198) (199) (200) (201) (202) (203) (204) (205) (206) (207) (208) (209) (210) (211) (212) (213) (214) (215) (216) (217) (218) (219) (220) (221) (222) (223) (224) (225) (226) (227) (228) (229) (230) (231) (232) (233) (234) (235) (236) (237) (238) (239) (240) (241) (242) (243) (244) (245) (246) (247) (248) (249) (250) (251) (252) (253) (254) (255) (256) (257) (258) (259) (260) (261) (262) (263) (264) (265) (266) (267) (268) (269) (270) (271) (272) (273) (274) (275) (276) (277) (278) (279) (280) (281) (282) (283) (284) (285) (286) (287) (288) (289) (290) (291) (292) (293) (294) (295) (296) (297) (298) (299) (300) (301) (302) (303) (304) (305) (306) (307) (308) (309) (310) (311) (312) (313) (314) (315) (316) (317) (318) (319) (320) (321) (322) (323) (324) (325) (326) (327) (328) (329) (330) (331) (332) (333) (334) (335) (336) (337) (338) (339) (340) (341) (342) (343) (344) (345) (346) (347) (348) (349) (350) (351) (352) (353) (354) (355) (356) (357) (358) (359) (360) (361) (362) (363) (364) (365) (366) (367) (368) (369) (370) (371) (372) (373) (374) (375) (376) (377) (378) (379) (380) (381) (382) (383) (384) (385) (386) (387) (388) (389) (390) (391) (392) (393) (394) (395) (396) (397) (398) (399) (400) (401) (402) (403) (404) (405) (406) (407) (408) (409) (410) (411) (412) (413) (414) (415) (416) (417) (418) (419) (420) (421) (422) (423) (424) (425) (426) (427) (428) (429) (430) (431) (432) (433) (434) (435) (436) (437) (438) (439) (440) (441) (442) (443) (444) (445) (446) (447) (448) (449) (450) (451) (452) (453) (454) (455) (456) (457) (458) (459) (460) (461) (462) (463) (464) (465) (466) (467) (468) (469) (470) (471) (472) (473) (474) (475) (476) (477) (478) (479) (480) (481) (482) (483) (484) (485) (486) (487) (488) (489) (490) (491) (492) (493) (494) (495) (496) (497) (498) (499) (500) (501) (502) (503) (504) (505) (506) (507) (508) (509) (510) (511) (512) (513) (514) (515) (516) (517) (518) (519) (520) (521) (522) (523) (524) (525) (526) (527) (528) (529) (530) (531) (532) (533) (534) (535) (536) (537) (538) (539) (540) (541) (542) (543) (544) (545) (546) (547) (548) (549) (550) (551) (552) (553) (554) (555) (556) (557) (558) (559) (560) (561) (562) (563) (564) (565) (566) (567) (568) (569) (570) (571) (572) (573) (574) (575) (576) (577) (578) (579) (580) (581) (582) (583) (584) (585) (586) (587) (588) (589) (590) (591) (592) (593) (594) (595) (596) (597) (598) (599) (600) (601) (602) (603) (604) (605) (606) (607) (608) (609) (610) (611) (612) (613) (614) (615) (616) (617) (618) (619) (620) (621) (622) (623) (624) (625) (626) (627) (628) (629) (630) (631) (632) (633) (634) (635) (636) (637) (638) (639) (640) (641) (642) (643) (644) (645) (646) (647) (648) (649) (650) (651) (652) (653) (654) (655) (656) (657) (658) (659) (660) (661) (662) (663) (664) (665) (666) (667) (668) (669) (670) (671) (672) (673) (674) (675) (676) (677) (678) (679) (680) (681) (682) (683) (684) (685) (686) (687) (688) (689) (690) (691) (692) (693) (694) (695) (696) (697) (698) (699) (700) (701) (702) (703) (704) (705) (706) (707) (708) (709) (710) (711) (712) (713) (714) (715) (716) (717) (718) (719) (720) (721) (722) (723) (724) (725) (726) (727) (728) (729) (730) (731) (732) (733) (734) (735) (736) (737) (738) (739) (740) (741) (742) (743) (744) (745) (746) (747) (748) (749) (750) (751) (752) (753) (754) (755) (756) (757) (758) (759) (760) (761) (762) (763) (764) (765) (766) (767) (768) (769) (770) (771) (772) (773) (774) (775) (776) (777) (778) (779) (780) (781) (782) (783) (784) (785) (786) (787) (788) (789) (790) (791) (792) (793) (794) (795) (796) (797) (798) (799) (800) (801) (802) (803) (804) (805) (806) (807) (808) (809) (810) (811) (812) (813) (814) (815) (816) (817) (818) (819) (820) (821) (822) (823) (824) (825) (826) (827) (828) (829) (830) (831) (832) (833) (834) (835) (836) (837) (838) (839) (840) (841) (842) (843) (844) (845) (846) (847) (848) (849) (850) (851) (852) (853) (854) (855) (856) (857) (858) (859) (860) (861) (862) (863) (864) (865) (866) (867) (868) (869) (870) (871) (872) (873) (874) (875) (876) (877) (878) (879) (880) (881) (882) (883) (884) (885) (886) (887) (888) (889) (890) (891) (892) (893) (894) (895) (896) (897) (898) (899) (900) (901) (902) (903) (904) (905) (906) (907) (908) (909) (910) (911) (912) (913) (914) (915) (916) (917) (918) (919) (920) (921) (922) (923) (924) (925) (926) (927) (928) (929) (930) (931) (932) (933) (934) (935) (936) (937) (938) (939) (940) (941) (942) (943) (944) (945) (946) (947) (948) (949) (950) (951) (952) (953) (954) (955) (956) (957) (958) (959) (960) (961) (962) (963) (964) (965) (966) (967) (968) (969) (970) (971) (972) (973) (974) (975) (976) (977) (978) (979) (980) (981) (982) (983) (984) (985) (986) (987) (988) (989) (990) (991) (992) (993) (994) (995) (996) (997) (998) (999) (1000).

• **Quantization-aware training (QAT) for [LoRa](#).** A recent study [QAT](#) explores the effect of [OTF](#) methods by applying a dual-bit distribution method to improve the weight activations as well as the bit-order scale by including post-train quantization based on [LoRa](#), for more precise results with the quantization on both weights and bit-order scales, but not on [LoRa](#) activation quantization, which still needs more experiments.

• **Efficient quantization can be directly applied to [LoRa](#).** While the performance of some previous weight quantization depends on specific methods (e.g., [OTF](#) (102) (103) (104) (105) (106) (107) (108) (109) (110) (111) (112) (113) (114) (115) (116) (117) (118) (119) (120) (121) (122) (123) (124) (125) (126) (127) (128) (129) (130) (131) (132) (133) (134) (135) (136) (137) (138) (139) (140) (141) (142) (143) (144) (145) (146) (147) (148) (149) (150) (151) (152) (153) (154) (155) (156) (157) (158) (159) (160) (161) (162) (163) (164) (165) (166) (167) (168) (169) (170) (171) (172) (173) (174) (175) (176) (177) (178) (179) (180) (181) (182) (183) (184) (185) (186) (187) (188) (189) (190) (191) (192) (193) (194) (195) (196) (197) (198) (199) (200) (201) (202) (203) (204) (205) (206) (207) (208) (209) (210) (211) (212) (213) (214) (215) (216) (217) (218) (219) (220) (221) (222) (223) (224) (225) (226) (227) (228) (229) (230) (231) (232) (233) (234) (235) (236) (237) (238) (239) (240) (241) (242) (243) (244) (245) (246) (247) (248) (249) (250) (251) (252) (253) (254) (255) (256) (257) (258) (259) (260) (261) (262) (263) (264) (265) (266) (267) (268) (269) (270) (271) (272) (273) (274) (275) (276) (277) (278) (279) (280) (281) (282) (283) (284) (285) (286) (287) (288) (289) (290) (291) (292) (293) (294) (295) (296) (297) (298) (299) (300) (301) (302) (303) (304) (305) (306) (307) (308) (309) (310) (311) (312) (313) (314) (315) (316) (317) (318) (319) (320) (321) (322) (323) (324) (325) (326) (327) (328) (329) (330) (331) (332) (333) (334) (335) (336) (337) (338) (339) (340) (341) (342) (343) (344) (345) (346) (347) (348) (349) (350) (351) (352) (353) (354) (355) (356) (357) (358) (359) (360) (361) (362) (363) (364) (365) (366) (367) (368) (369) (370) (371) (372) (373) (374) (375) (376) (377) (378) (379) (380) (381) (382) (383) (384) (385) (386) (387) (388) (389) (390) (391) (392) (393) (394) (395) (396) (397) (398) (399) (400) (401) (402) (403) (404) (405) (406) (407) (408) (409) (410) (411) (412) (413) (414) (415) (416) (417) (418) (419) (420) (421) (422) (423) (424) (425) (426) (427) (428) (429) (430) (431) (432) (433) (434) (435) (436) (437) (438) (439) (440) (441) (442) (443) (444) (445) (446) (447) (448) (449) (450) (451) (452) (453) (454) (455) (456) (457) (458) (459) (460) (461) (462) (463) (464) (465) (

- [illegible]