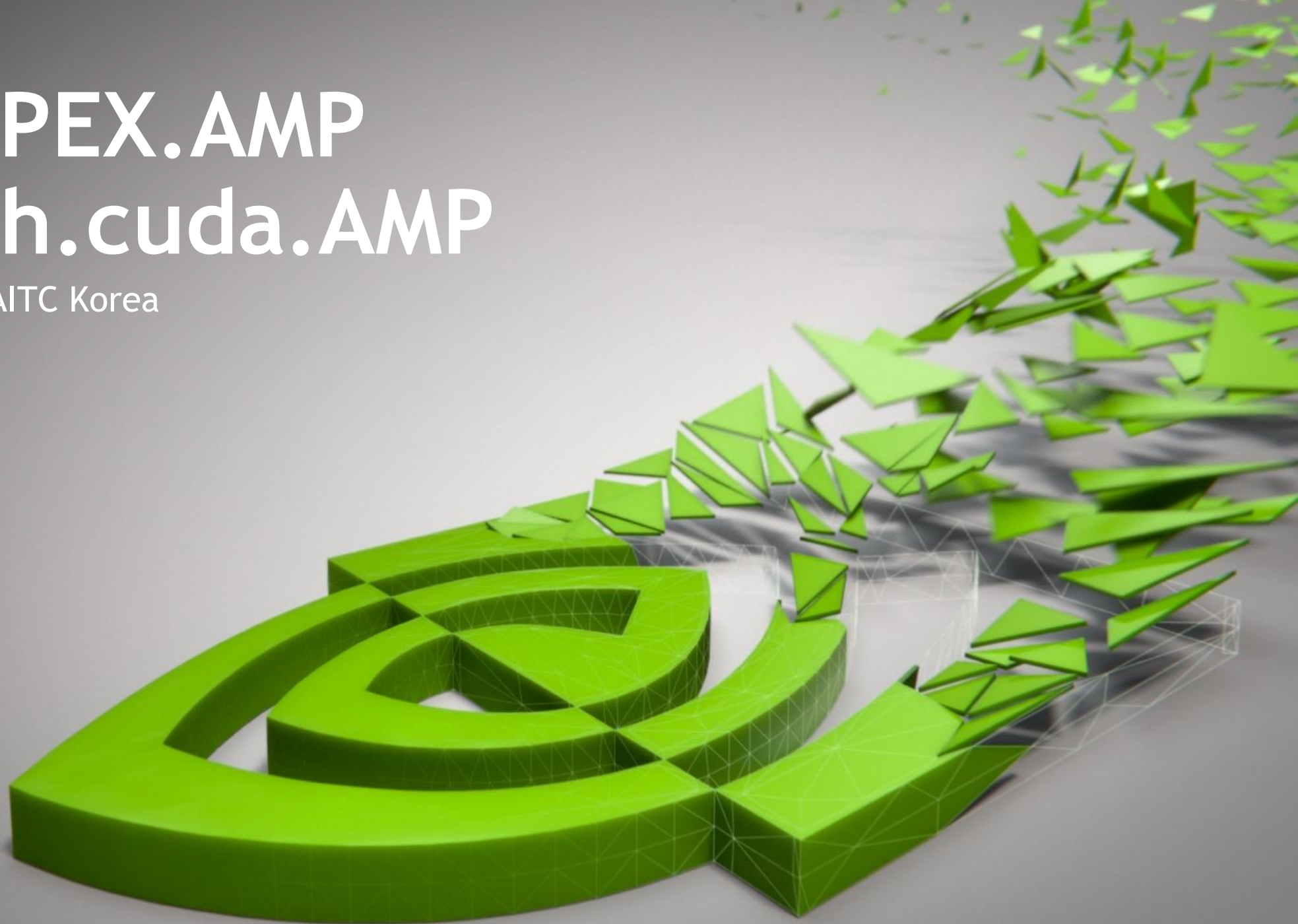# From APEX.AMP
# To torch.cuda.AMP

Hyungon Ryu | NVAITC Korea

# AGENDA
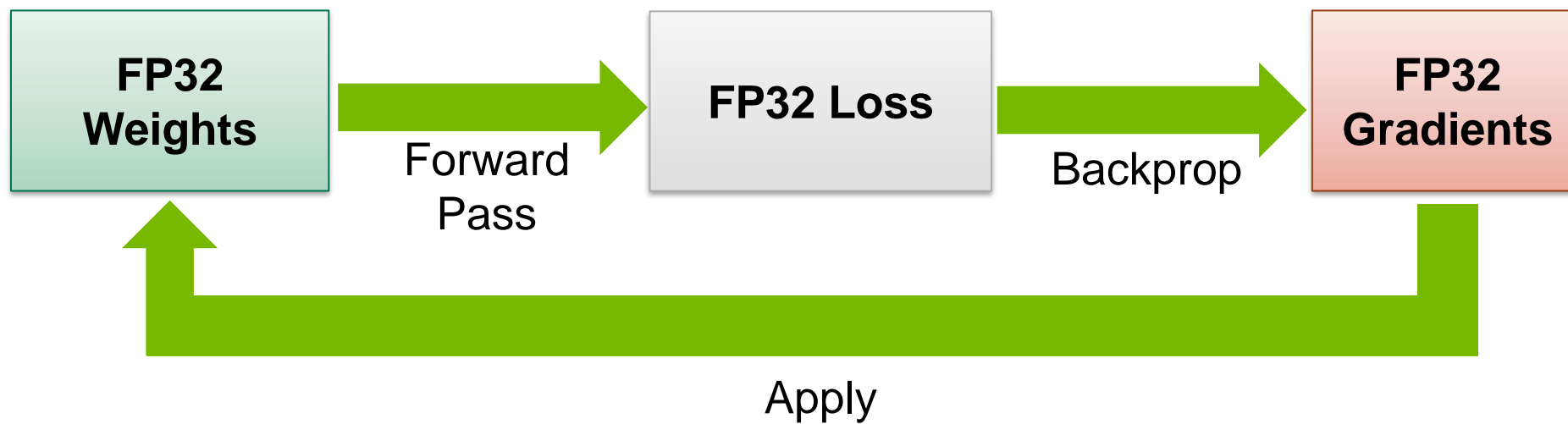
**Precision**

**apex.amp**
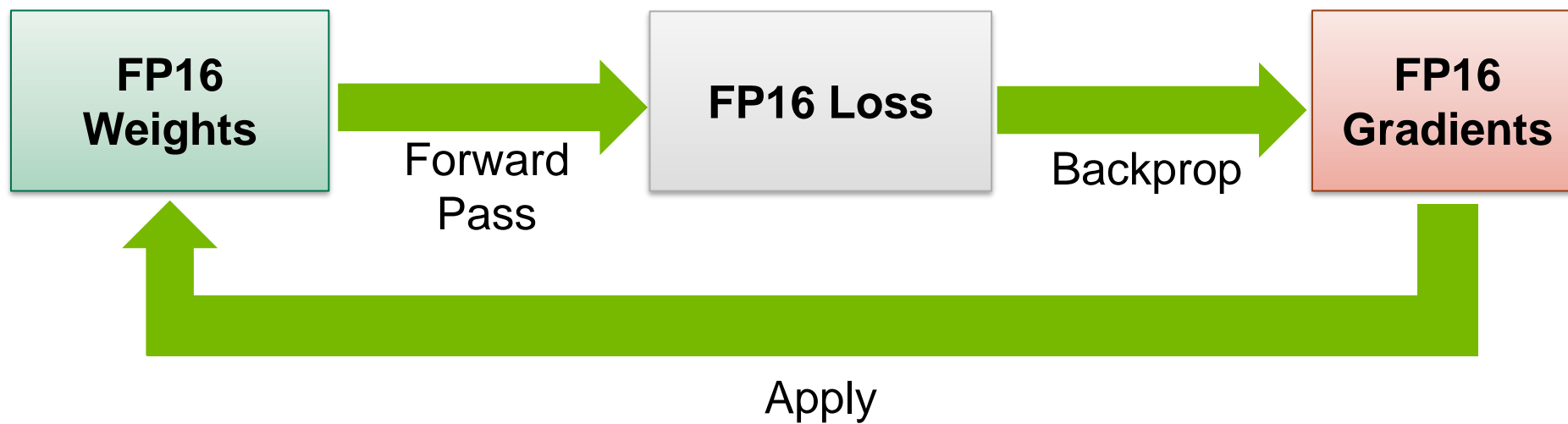
**torch.cuda.amp**

NVIDIA.

# PRECISION

# FP32 TRAINING

# FP16 TRAINING

# PRECISION FOR DL



Sign | Range | Precision

**FP32** — 8 BITS | 23 BITS — **32Bit**

TF32 Range

**TENSOR FLOAT 32 (TF32)** — 8 BITS | 10 BITS — **19Bit**

TF32 Precision

**FP16** — 5 BITS | 10 BITS

**16Bit**

**BFLOAT16** — 8 BITS | 7 BITS
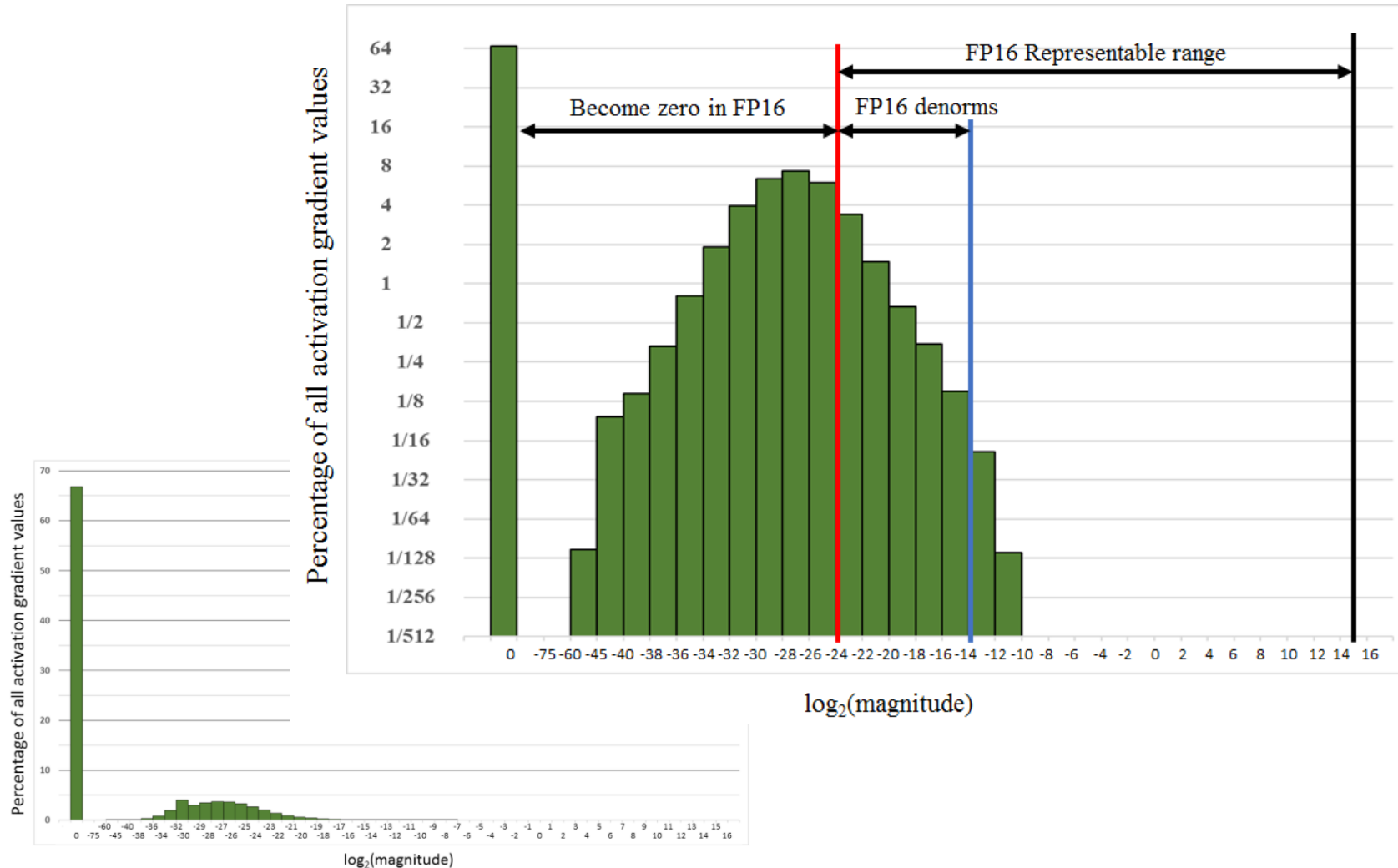
**Benefit of FP16**
Lowering required memory
~ large mini-batches
~ less data transfer
Faster training/inferencing

# A HISTOGRAM OF ACTIVATION GRADIENTS

## Representation of Half Precision floating point format

# GRADIENTS MAY UNDERFLOW

# GRADIENT EXPLODE



Legend:
- FP32
- Mixed Precision, loss scale 1
- Mixed Precision, loss scale 128

FP16
Mixed Precision

X-axis: Training Iteration (0K, 500K, 1,000K, 1,500K, 2,000K)
Y-axis: 2.5, 3.0, 3.5, 4.0, 4.5, 5.0

# MIXED PRECISION TRAINING

**Compute in FP16**
**Save in FP32**

# MIXED PRECISION TRAINING

## Compute in FP16
## Save in FP32

APEX.AMP

# APEX.AMP

https://nvidia.github.io/apex/amp.html

APEX enable DDP, Mixed Precision(FP16) and AMP

**amp.initialize(model, optimizer, opt_level="O1")**

        O0 : FP32, O1/O2 : AMP,  O3 : FP16

with amp.scale_loss(loss, optimizer) as scaled_loss:

    scaled_loss.backward()

# 3 STEP FOR APEX.AMP

Load APEX.AMP module

Initiate APEX AMP

Enable Loss scale

# APEX.AMP
## Step1. load APEX module

```python
39    import torch.distributed as dist
40    from torch.utils.data.distributed import DistributedSampler
41
42    from apex.parallel import DistributedDataParallel as DDP
43
44    import models
45    import loss_functions
46    import data_functions
47    from tacotron2_common.utils import ParseFromConfigFile
48
49    import dllogger as DLLogger
50    from dllogger import StdOutBackend, JSONStreamBackend, Verbosity
51
52    from scipy.io.wavfile import write as write_wav
53
54    from apex import amp
55    amp.lists.functional_overrides.FP32_FUNCS.remove('softmax')
56    amp.lists.functional_overrides.FP16_FUNCS.append('softmax')
57
```

# APEX.AMP
## Step2. initiate AMP

```
387    model_config = models.get_model_config(model_name, args)
388    model = models.get_model(model_name, model_config,
389                             cpu_run=False,
390                             uniform_initialize_bn_weight=not args.disable_uniform_init:
391
392    if not args.amp and distributed_run:
393        model = DDP(model)
394
395    optimizer = torch.optim.Adam(model.parameters(), lr=args.learning_rate,
396                                 weight_decay=args.weight_decay)
397
398    if args.amp:
399        model, optimizer = amp.initialize(model, optimizer, opt_level="O1")
400        if distributed_run:
401            model = DDP(model)
402
```

# APEX.AMP
## Step3. enable loss scale

```
481        x, y, num_items = batch_to_gpu(batch)
482
483        y_pred = model(x)
484        loss = criterion(y_pred, y)
485
486        if distributed_run:
487            reduced_loss = reduce_tensor(loss.data, world_size).item()
488            reduced_num_items = reduce_tensor(num_items.data, 1).item()
489        else:
490            reduced_loss = loss.item()
491            reduced_num_items = num_items.item()



502        if args.amp:
503            with amp.scale_loss(loss, optimizer) as scaled_loss:
504                scaled_loss.backward()
505            grad_norm = torch.nn.utils.clip_grad_norm_(
506                amp.master_params(optimizer), args.grad_clip_thresh)
507        else:
508            loss.backward()
509            grad_norm = torch.nn.utils.clip_grad_norm_(
510                model.parameters(), args.grad_clip_thresh)
511
512        optimizer.step()
```

# ALSO NEED

## amp.state_dict()

```python
191    def save_checkpoint(model, optimizer, epoch, config, amp_run, output_dir, model_name,
192                        local_rank, world_size):
193
194        random_rng_state = torch.random.get_rng_state().cuda()
195        cuda_rng_state = torch.cuda.get_rng_state(local_rank).cuda()
196
197        random_rng_states_all = [torch.empty_like(random_rng_state) for _ in range(world_size)]
198        cuda_rng_states_all = [torch.empty_like(cuda_rng_state) for _ in range(world_size)]
199
200        if world_size > 1:
201            dist.all_gather(random_rng_states_all, random_rng_state)
202            dist.all_gather(cuda_rng_states_all, cuda_rng_state)
203        else:
204            random_rng_states_all = [random_rng_state]
205            cuda_rng_states_all = [cuda_rng_state]
206
207        random_rng_states_all = torch.stack(random_rng_states_all).cpu()
208        cuda_rng_states_all = torch.stack(cuda_rng_states_all).cpu()
209
210        if local_rank == 0:
211            checkpoint = {'epoch': epoch,
212                          'cuda_rng_state_all': cuda_rng_states_all,
213                          'random_rng_states_all': random_rng_states_all,
214                          'config': config,
215                          'state_dict': model.state_dict(),
216                          'optimizer': optimizer.state_dict()}
217            if amp_run:
218                checkpoint['amp'] = amp.state_dict()
219
```

```python
246    def load_checkpoint(model, optimizer, epoch, config, amp_run, filepath, local_rank):
247
248        checkpoint = torch.load(filepath, map_location='cpu')
249
250        epoch[0] = checkpoint['epoch']+1
251        device_id = local_rank % torch.cuda.device_count()
252        torch.cuda.set_rng_state(checkpoint['cuda_rng_state_all'][device_id])
253        if 'random_rng_states_all' in checkpoint:
254            torch.random.set_rng_state(checkpoint['random_rng_states_all'][device_id])
255        elif 'random_rng_state' in checkpoint:
256            torch.random.set_rng_state(checkpoint['random_rng_state'])
257        else:
258            raise Exception("Model checkpoint must have either 'random_rng_state' or 'random_rng_states_all' key.")
259        config = checkpoint['config']
260        model.load_state_dict(checkpoint['state_dict'])
261        optimizer.load_state_dict(checkpoint['optimizer'])
262
263        if amp_run:
264            amp.load_state_dict(checkpoint['amp'])
265
```

# TORCH.CUDA.AMP

# TORCH NATIVE AMP
## https://pytorch.org/docs/stable/amp.html

APEX.AMP Upstream  to pytorch ( Oct. 2020)


torch.cuda.amp.autocast(*enabled=True*)

torch.cuda.amp.GradScaler(*init_scale=65536.0, growth_factor=2.0, backoff_factor=0.5, growth_interval=2000, enabled=True*)

# TORCH.CUDA.AMP
## Step1. enable scaler

```python
871        model = model.to(device)
872
873        scaler = None
874        if args.fp16:
875            if args.amp == 'pytorch':
876                scaler = torch.cuda.amp.GradScaler()
877            elif args.amp == 'apex':
878                model, optimizer = amp.initialize(
879                    model,
880                    optimizer,
881                    opt_level=args.apex_amp_opt_level,
882                    )
883
884        if args.multi_gpu == 'ddp' and torch.distributed.is_initialized():
885            para_model = DistributedDataParallel(model,
886                                                 device_ids=[args.local_rank],
887                                                 output_device=args.local_rank,
888                                                 broadcast_buffers=False,
889                                                 find_unused_parameters=True,
```

# TORCH.CUDA.AMP
## Step2. enable autocast

```python
def train_iteration(model, i, mems, data_chunks, target_chunks, scaler,
                    optimizer, device, delay_unscale, args):
    cpu = torch.device('cpu')
    data_i = data_chunks[i].contiguous()
    target_i = target_chunks[i].contiguous()

    if args.swap_mem and mems[i] is not None:
        mems[i] = mems[i].to(device, non_blocking=True)

    enable_autocast = args.fp16 and args.amp == 'pytorch'
    with torch.cuda.amp.autocast(enable_autocast):
        loss, mems[i] = model(data_i, target_i, mems[i])
        loss = loss.float().mean().type_as(loss) / args.batch_chunk

    if args.swap_mem and mems[i] is not None:
        mems[i] = mems[i].to(cpu, non_blocking=True)

    if args.fp16:
        if args.amp == 'pytorch':
            scaler.scale(loss).backward()
        elif args.amp == 'apex':
            with amp.scale_loss(loss, optimizer, delay_unscale=delay_unscale) as scaled_loss:
                scaled_loss.backward()
    else:
        loss.backward()

    train_loss = loss.float().item()
    return train_loss
```

# TORCH.CUDA.AMP
## Step3. use scaler

```python
464    def train_iteration(model, i, mems, data_chunks, target_chunks, scaler,
465                        optimizer, device, delay_unscale, args):
466        cpu = torch.device('cpu')
467        data_i = data_chunks[i].contiguous()
468        target_i = target_chunks[i].contiguous()
469
470        if args.swap_mem and mems[i] is not None:
471            mems[i] = mems[i].to(device, non_blocking=True)
472
473        enable_autocast = args.fp16 and args.amp == 'pytorch'
474        with torch.cuda.amp.autocast(enable_autocast):
475            loss, mems[i] = model(data_i, target_i, mems[i])
476            loss = loss.float().mean().type_as(loss) / args.batch_chunk
477
478        if args.swap_mem and mems[i] is not None:
479            mems[i] = mems[i].to(cpu, non_blocking=True)
480
481        if args.fp16:
482            if args.amp == 'pytorch':
483                scaler.scale(loss).backward()
484            elif args.amp == 'apex':
485                with amp.scale_loss(loss, optimizer, delay_unscale=delay_unscale) as scaled_loss:
486                    scaled_loss.backward()
487        else:
488            loss.backward()
489
490        train_loss = loss.float().item()
491        return train_loss
492
```

# ALSO NEED

## scaler.state_dict()

```python
296  def save_checkpoint(args, model, model_config, optimizer, scheduler, scaler,
297                      vocab, epoch, batch, last_iter, train_step, best_val_loss,
298                      is_best, work_dir):
299      if args.fp16:
300          if args.amp == 'pytorch':
301              amp_state = scaler.state_dict()
302          elif args.amp == 'apex':
303              amp_state = amp.state_dict()
304      else:
305          amp_state = None
306
307      state = {
308          'args': args,
309          'model_config': model_config,
310          'model_state': model.state_dict(),
311          'optimizer_state': optimizer.state_dict(),
312          'scheduler_state': scheduler.state_dict(),
313          'vocab': vocab,
314          'amp_state': amp_state,
315          'epoch': epoch,
316          'batch': batch,
317          'last_iter': last_iter,
318          'train_step': train_step,
319          'best_val_loss': best_val_loss,
320          }
```

# ADDITIONAL TIPS

Do not check loss.item() for every iteration