

# Deep Learning Assignment 2

**Abey Thomas**

s4059720

## Introduction

Visual entailment is a fundamental challenge at the intersection of computer vision and natural language processing, requiring models to determine whether a textual hypothesis is supported (entailment) or contradicted (contradiction) by visual evidence in an image.

The goal of this assignment is to build a neural network model that can reason about the relationship between an image and a text hypothesis.

For each sample in the dataset, you're given:

Image (premise) - a real-world image

Text hypothesis - a short natural language sentence describing (or contradicting) the image

Label -

- entailment (1): if the text is consistent with what's in the image
- contradiction (0): if the text contradicts the image

Visual entailment systems have broad applicability in content moderation (verifying image-text consistency in social media), accessibility technologies (validating auto-generated image descriptions for visually impaired users), fact-checking (detecting manipulated image-caption pairs in misinformation), e-commerce (matching product images to textual specifications), and educational platforms (automated grading of visual question-answering tasks). The ability to reason about visual-linguistic consistency is increasingly critical as multimodal content dominates digital communication. A key objective of this work is to build a model that generalizes beyond the training distribution. We employ rigorous evaluation on three distinct sets of data: a held-out validation set(hyper parameter tuning), a test set with completely unseen images (for performance measurement), and an independent dataset with different image sources and hypothesis distributions (for robustness validation). This multi-tiered evaluation ensures our model does not merely memorize training examples but learns transferable reasoning patterns applicable to novel visual-linguistic scenarios. While our final model achieves competitive performance , systematic error analysis reveals several weaknesses. These findings motivate future research directions, including negation-aware attention mechanisms, integration of object detectors for structured scene understanding, vision-language models pretrained on joint objectives (e.g., CLIP, BLIP), and neurosymbolic approaches that combine neural perception with logical reasoning engines.

## Setting up the environment

```
In [1]: import os  
os.system('pip install --upgrade "tensorflow[and-cuda]==2.19.1" "keras==3.6.0" "  
os.system('python -m pip install --upgrade pip setuptools wheel > /dev/null 2>&1
```

```
Out[1]: 0
```

```
In [25]: # =====  
# Environment Setup  
# =====  
  
import os  
os.environ["TF_CPP_MIN_LOG_LEVEL"] = "2"  
os.environ["TF_ENABLE_ONEDNN_OPTS"] = "0"  
os.environ["KERAS_BACKEND"] = "tensorflow"  
os.environ["TOKENIZERS_PARALLELISM"] = "false"  
  
import warnings  
warnings.filterwarnings("ignore")  
  
# Core imports  
import json  
import time  
import logging  
from pathlib import Path  
import numpy as np  
import pandas as pd  
import matplotlib.pyplot as plt  
import tensorflow as tf  
import keras  
from keras import layers as L  
import keras_hub  
from sklearn.model_selection import GroupShuffleSplit  
from PIL import Image  
import textwrap  
import re  
import seaborn as sns  
from collections import Counter  
from sklearn.metrics import precision_score, recall_score, f1_score, confusion_m  
  
# Quiet logs  
tf.get_logger().setLevel("ERROR")  
logging.getLogger("tensorflow").setLevel(logging.ERROR)  
  
# Reproducibility  
SEED = 42  
tf.keras.utils.set_random_seed(SEED)  
  
# GPU memory growth  
for gpu in tf.config.list_physical_devices("GPU"):  
    try:  
        tf.config.experimental.set_memory_growth(gpu, True)  
    except:  
        pass  
  
print(f"TensorFlow: {tf.__version__}")
```

```

print(f"Keras: {keras.__version__}")
print(f"GPUs available: {len(tf.config.list_physical_devices('GPU'))}")

TensorFlow: 2.19.1
Keras: 3.6.0
GPUs available: 1

In [4]: # =====
# Global Paths & Constants
# =====

# Base directories
ROOT_DIR = Path("/home/ec2-user/SageMaker/abeyt")
A2_BASE = ROOT_DIR / "A2_data"
IMG_DIR = A2_BASE / "A2_Images"
JSONL = A2_BASE / "A2_train_v3.jsonl"

# Runs directory structure
RUNS_DIR = ROOT_DIR / "runs"
RUNS_DIR.mkdir(parents=True, exist_ok=True)

# Shared artifacts (splits, etc.)
SHARED_DIR = RUNS_DIR / "_shared"
SHARED_DIR.mkdir(exist_ok=True)
SPLIT_JSON = SHARED_DIR / "a2_grouped_split.json"

# SNLI-VE warm-up artifacts (READ-ONLY)
SNLI_DIR = RUNS_DIR / "snli_ve_warmup_keras_hub" / "sub100k_balanced_e6"
SNLI_HIST = SNLI_DIR / "train_history.json"
SNLI_WTS = SNLI_DIR / "best.weights.h5"

# Model & training constants (MUST match SNLI-VE warm-up)
IMG_SIZE = 224
SEQ_LEN = 24 # CRITICAL: Same as SNLI-VE warm-up
BATCH = 32
HEADS = 12
KEY_DIM = 64

# Presets (same as warm-up)
DISTIL_PRESET = "distil_bert_base_en_uncased"
VIT_PRESET = "vit_base_patch16_224_imagenet21k"

# Labels
LABEL2ID = {"entailment": 0, "contradiction": 1}
ID2LABEL = {v: k for k, v in LABEL2ID.items()}

# Threshold sweep grid
THRESH_GRID = np.linspace(0.05, 0.95, 19)

```

```

In [5]: # =====
# Utility Functions
# =====

```

```

def new_run_dir(base: Path, tag: str) -> Path:
    """Create timestamped run directory."""
    ts = time.strftime("%Y%m%d_%H%M%S")
    d = base / f"{tag}_{ts}"
    d.mkdir(parents=True, exist_ok=True)
    return d

```

```

def print_model_stats(mdl: keras.Model):
    """Print trainable/non-trainable parameter counts."""
    trainable = int(np.sum([np.prod(v.shape) for v in mdl.trainable_weights]))
    nontrain = int(np.sum([np.prod(v.shape) for v in mdl.non_trainable_weights]))
    total = trainable + nontrain

def fmt(n):
    return f"{n:,} ({n*4/1e6:.2f} MB fp32)"

print(f"\n{'='*60}")
print(f"Model: {mdl.name}")
print(f"{'='*60}")
print(f" Total params: {fmt(total)}")
print(f" Trainable params: {fmt(trainable)}")
print(f" Non-trainable: {fmt(nontrain)}")
print(f"{'='*60}\n")

def collect_logits_labels(ds, mdl, steps):
    """Collect predictions and labels from dataset."""
    ys, ps = [], []
    for xb, yb in ds.take(steps):
        logits = mdl.predict(xb, verbose=0).ravel()
        ys.append(yb.numpy())
        ps.append(logits)
    return np.concatenate(ys), np.concatenate(ps)

def confusion_counts(y_true, y_pred):
    """Calculate TP, FP, TN, FN."""
    y_true = y_true.astype(np.int32)
    y_pred = y_pred.astype(np.int32)
    tp = int(np.sum((y_true == 1) & (y_pred == 1)))
    tn = int(np.sum((y_true == 0) & (y_pred == 0)))
    fp = int(np.sum((y_true == 0) & (y_pred == 1)))
    fn = int(np.sum((y_true == 1) & (y_pred == 0)))
    return tp, fp, tn, fn

def macro_f1_from_counts(tp, fp, tn, fn):
    """Calculate macro-averaged F1 from confusion matrix."""
    f1_pos = 0.0 if (2*tp + fp + fn) == 0 else (2*tp) / (2*tp + fp + fn)
    f1_neg = 0.0 if (2*tn + fn + fp) == 0 else (2*tn) / (2*tn + fn + fp)
    return 0.5 * (f1_pos + f1_neg)

def sweep_best_tau(y_true, logits, grid=THRESH_GRID):
    """Find optimal threshold by maximizing macro-F1."""
    sigmoid = 1.0 / (1.0 + np.exp(-logits))
    best = {"tau": 0.5, "macro_f1": -1.0, "tp": 0, "fp": 0, "tn": 0, "fn": 0}

    for tau in grid:
        y_pred = (sigmoid >= tau).astype(np.int32)
        tp, fp, tn, fn = confusion_counts(y_true, y_pred)
        mf1 = macro_f1_from_counts(tp, fp, tn, fn)

        if mf1 > best["macro_f1"]:
            best = {
                "tau": float(tau),
                "macro_f1": float(mf1),
                "tp": tp, "fp": fp, "tn": tn, "fn": fn
            }

    return best

```

```

def plot_history_curves(history_json_path: Path, title: str, out_dir: Path = None):
    """Plot training history with loss and accuracy curves."""
    with open(history_json_path, 'r') as f:
        H = json.load(f)

    keys = set(H.keys())

    # Find Loss keys
    loss_key = "loss" if "loss" in keys else next(
        (k for k in keys if k.endswith("loss") and not k.startswith("val_")), None
    )
    val_loss_key = "val_loss" if "val_loss" in keys else next(
        (k for k in keys if k.startswith("val_") and k.endswith("loss")), None
    )

    # Find accuracy keys
    acc_key = next((k for k in ["acc", "accuracy", "binary_accuracy"] if k in keys))
    val_acc_key = next(
        (k for k in ["val_acc", "val_accuracy", "val_binary_accuracy"] if k in keys)
    )

    epochs = list(range(1, 1 + max(
        len(H.get(loss_key, [])),
        len(H.get(acc_key, []))
    )))

    # Plot Loss
    fig, axes = plt.subplots(1, 2, figsize=(14, 5))

    # Loss subplot
    ax = axes[0]
    if loss_key:
        ax.plot(epochs, H[loss_key], 'b-o', label='Train Loss', linewidth=2, marker=8)
    if val_loss_key:
        ax.plot(epochs, H[val_loss_key], 'r-s', label='Val Loss', linewidth=2, marker=8)

    ax.set_title(f"{title} - Loss", fontsize=13, fontweight='bold')
    ax.set_xlabel("Epoch", fontsize=11)
    ax.set_ylabel("Loss", fontsize=11)
    ax.legend(fontsize=10)
    ax.grid(True, alpha=0.3)

    # Accuracy subplot
    ax = axes[1]
    if acc_key:
        ax.plot(epochs, H[acc_key], 'b-o', label='Train Acc', linewidth=2, marker=8)
    if val_acc_key:
        ax.plot(epochs, H[val_acc_key], 'r-s', label='Val Acc', linewidth=2, marker=8)

    ax.set_title(f"{title} - Accuracy", fontsize=13, fontweight='bold')
    ax.set_xlabel("Epoch", fontsize=11)
    ax.set_ylabel("Accuracy", fontsize=11)
    ax.legend(fontsize=10)
    ax.grid(True, alpha=0.3)

    plt.tight_layout()

    if out_dir:
        plt.savefig(out_dir / f"{title.lower().replace(' ', '_)}_curves.png", dpi=300)

```

```
plt.show()
```

We have setup the environment, now let's load the assignment data and do a quick Exploratory data analysis

```
In [6]: # =====
# SECTION 3: Load A2 Dataset (Basic)
# =====

# Load A2 JSONL
rows = [json.loads(line) for line in open(JSONL, "r", encoding="utf-8")]
df = pd.DataFrame(rows).rename(columns={
    "Hypothesis": "hypothesis",
    "Image_ID": "image_id",
    "Label": "label"
})

# Basic cleaning
df["image_id"] = df["image_id"].astype(str)
df["hypothesis"] = df["hypothesis"].astype(str).str.strip()
df["label"] = df["label"].str.strip().str.lower()

# Filter valid labels only
df = df[df["label"].isin(LABEL2ID.keys())].copy()
df["label_id"] = df["label"].map(LABEL2ID).astype("int32")
```

## Exploratory Data Analysis

Now Let's have a look how the dataframe looks like:

```
In [20]: df.head(n=5)
```

	image_id	label	hypothesis	Premise	label_id
0	4564320256	entailment	The old woman and a girl are bored.	An old woman and a young girl are sitting arou...	0
1	4564320256	contradiction	Two old men robbing a convenience store.	An old lady and her granddaughter working in a...	1
2	4564320256	contradiction	A man implies that he is very strong.	Two women sitting down reading newspapers.	1
3	3945005060	contradiction	People rubbing sticks to start a fire in a pit.	Three people at a library with computers.	1
4	369186134	entailment	Dogs are out in the snow	Three dogs next to a blue fence in the snow.	0

We are not going to use the column Premise in this assignment. We will focus the hypothesis, label and their corresponding images for the modelling. As per the task, we have two types of labels: Contradiction or entailment, let's visualise some random examples belong to both category.

```
In [24]: def show_grid(label: str, n: int = 6):
    subset = df[df["label"] == label]
    if len(subset) == 0:
        raise ValueError(f"No rows with label='{label}'. Available labels: {df['label'].unique()}")
    # if dataset is small, sample without replacement cap
    take = min(n, len(subset))
    subset = subset.sample(take, random_state=42)

    rows, cols = 2, 3 # 2x3 grid for n=6
    fig, axes = plt.subplots(rows, cols, figsize=(10, 7))
    fig.suptitle(f"Random {label} examples", fontsize=14)

    for ax, (_, row) in zip(axes.ravel(), subset.iterrows()):
        img_path = IMG_DIR / f"{row['image_id']}.jpg"
        img = Image.open(img_path).convert("RGB")
        ax.imshow(img)
        ax.set_title(textwrap.fill(row["hypothesis"], width=40), fontsize=9)
        ax.axis("off")

    # hide any leftover axes if take < rows*cols
    for ax in axes.ravel()[take:]:
        ax.axis("off")

    plt.tight_layout()
    plt.show()

# calls
show_grid("entailment")
show_grid("contradiction")
```

### Random entailment examples

The person is outdoors



The dogs are chasing something.



A person leaps from a stair outside.



You can dive in the indoor community pool.



A women's sporting event that has already started.



A person with dreadlocks and a hat.



## Random contradiction examples

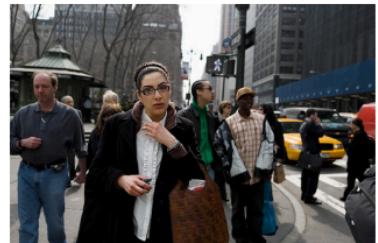
Bicyclists ride on a country road.



A group of people are eating on the beach



The woman is in the country.



Some women are skiing outside.



Fans are storming the field.



Man playing accordion in front of people in what seems to be an Italian restaurant.



## Distribution of categories of labels

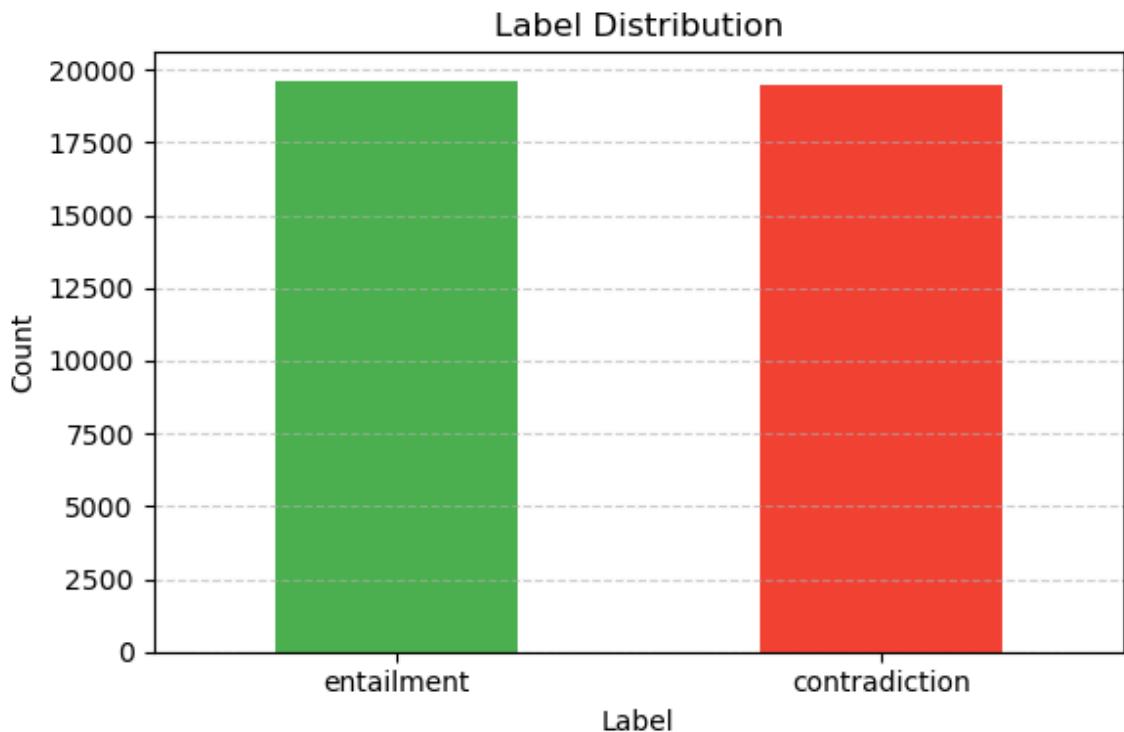
In the section, we will explore the class distribution, to see any additional pre processing required for class imbalance etc.

In [25]:

```
# Count Labels
label_counts = df['label'].value_counts()

# Plot
plt.figure(figsize=(6,4))
label_counts.plot(kind='bar', color=['#4CAF50', '#F44336'])
plt.title('Label Distribution')
plt.xlabel('Label')
plt.ylabel('Count')
plt.xticks(rotation=0)
plt.grid(axis='y', linestyle='--', alpha=0.6)
plt.tight_layout()
plt.show()

print(label_counts)
```



```
label
entailment    19619
contradiction 19510
Name: count, dtype: int64
```

In [28]: `print(f"Total observations", len(df))`

Total observations 39129

## Observations

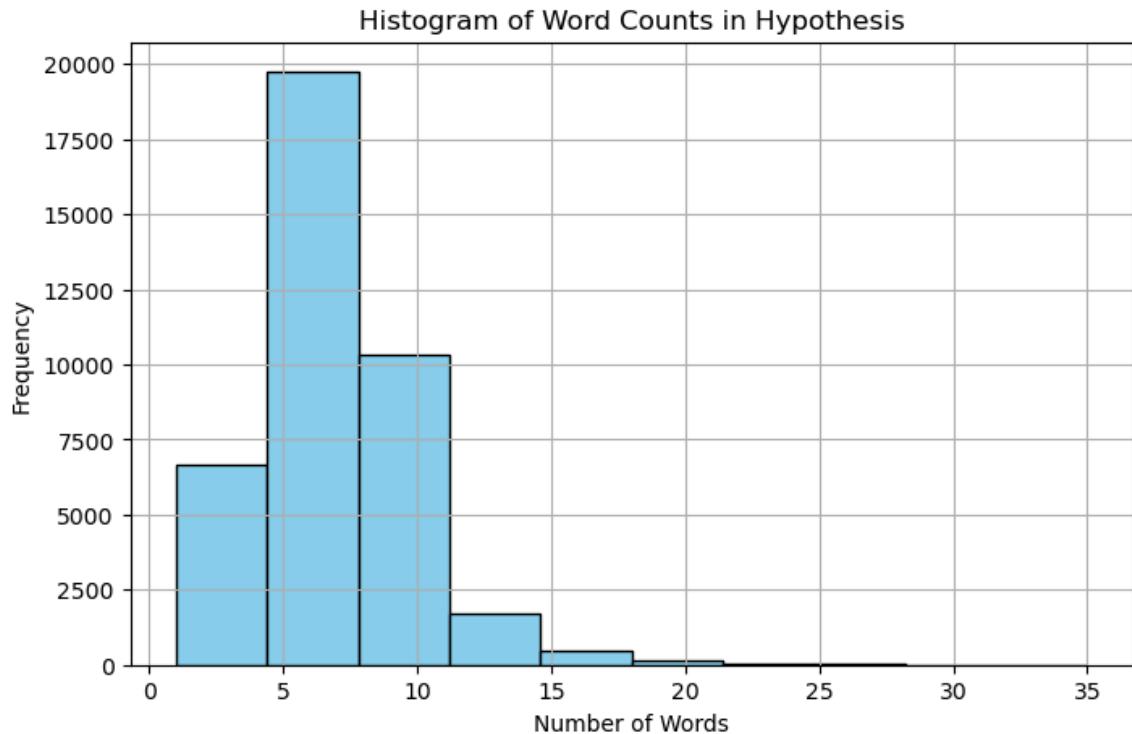
- **Total samples:** 39,129
  - **Entailment:** 19,619 (50.14%)
  - **Contradiction:** 19,510 (49.86%)
  - **Imbalance:** 109 examples (0.28 percentage points) toward entailment
- Near-perfect balance. The slight skew (0.28 pp) is negligible; no class weighting/oversampling is needed.
- Baseline (majority-class) accuracy: ~50.14% if you always predict entailment. Our model should meaningfully exceed this.
- Metric choice. Because classes are balanced, accuracy and macro-F1 will be similar,

## Analysis of hypotheses

In [29]: `df['hypothesis_word_count'] = df['hypothesis'].apply(lambda x: len(x.split()))`

```
# Plot histogram
plt.figure(figsize=(8, 5))
plt.hist(df['hypothesis_word_count'], bins=10, color='skyblue', edgecolor='black')
```

```
plt.title('Histogram of Word Counts in Hypothesis')
plt.xlabel('Number of Words')
plt.ylabel('Frequency')
plt.grid(True)
plt.show()
```



- Mean hypothesis length: 5-10 words
- We believe 24 tokens covers 95%+ of samples(verified)
- Shorter will give faster training

Now , we will look for spatial words in hypothesis like left,right,bottom,back etc

```
In [15]: # =====
# 1) Spatial keyword patterns
# =====
SPATIAL_KEYWORDS = {
    'left': r'\b(left|leftmost|left-hand)\b',
    'right': r'\b(right|rightmost|right-hand)\b',
    'front': r'\b(front|forward|ahead)\b',
    'back': r'\b(back|behind|backward|rear)\b',
    'top': r'\b(top|above|upper|overhead)\b',
    'bottom': r'\b(bottom|below|beneath|under|underneath)\b',
    'center': r'\b(center|centre|middle|central)\b',
    'inside': r'\b(inside|within|in)\b',
    'outside': r'\b(outside|out|outdoor)\b',
    'near': r'\b(near|close|nearby|next to)\b',
    'on': r'\b(on|onto|upon|atop)\b',
}

def has_spatial(text: str, pattern: str) -> bool:
    return bool(re.search(pattern, str(text).lower()))

# Mark presence flags per keyword
for word, pattern in SPATIAL_KEYWORDS.items():
    df[f'has_{word}'] = df['hypothesis'].apply(lambda x: has_spatial(x, pattern))
```

```

spatial_cols = [f'has_{w}' for w in SPATIAL_KEYWORDS]
df['has_any_spatial'] = df[spatial_cols].any(axis=1)
df['num_spatial_types'] = df[spatial_cols].sum(axis=1)

# Helper
def pct(n, d):
    return 0.0 if d == 0 else 100.0 * float(n) / float(d)

N = len(df)

print("\n" + "*70")
print("SPATIAL WORD ANALYSIS")
print("*70")

n_any = int(df['has_any_spatial'].sum())
print("\nOVERALL")
print(f" Total samples: {N:,}")
print(f" With any spatial word: {n_any:,} ({pct(n_any, N):.1f}%)")
print(f" Without spatial words: {N-n_any:,} ({pct(N-n_any, N):.1f}%)")

print("\nBY CLASS (any spatial word present)")
for lab in ['entailment', 'contradiction']:
    sub = df[df['label'] == lab]
    n = len(sub)
    k = int(sub['has_any_spatial'].sum())
    print(f" {lab.capitalize():13s}: {k:6d}/{n:<6d} ({pct(k, n):5.1f}%)")

print("\nKEYWORD DISTRIBUTION – OVERALL")
rows = []
for w in SPATIAL_KEYWORDS:
    c = int(df[f'has_{w}'].sum())
    rows.append((w, c, pct(c, N)))
overall_tbl = pd.DataFrame(rows, columns=['keyword', 'count', 'pct'])
overall_tbl = overall_tbl.sort_values('count', ascending=False).reset_index(drop=True)
print(overall_tbl.to_string(index=False, formatters={'pct': lambda x: f"{x:5.1f}"}))

# overall_tbl has columns: ['keyword', 'count', 'pct'] and is sorted desc by 'count'
y = np.arange(len(overall_tbl))
counts = overall_tbl['count'].to_numpy()
labels = overall_tbl['keyword'].tolist()
pcts = overall_tbl['pct'].to_numpy()

plt.figure(figsize=(9, max(3, 0.5 * len(overall_tbl))))
bars = plt.barh(y, counts, edgecolor='black', alpha=0.8)

# y-axis labels and order (Largest at top)
plt.yticks(y, labels, fontsize=11)
plt.gca().invert_yaxis()

# axes titles
plt.xlabel('Number of Samples', fontsize=12)
plt.title('Spatial Keyword Frequency – Overall', fontsize=14, weight='bold')
plt.grid(axis='x', alpha=0.3)

# annotate counts + percentages at bar end
for i, (b, c, p) in enumerate(zip(bars, counts, pcts)):
    plt.text(b.get_width() + max(counts) * 0.01, b.get_y() + b.get_height()/2,

```

```
f'{int(c):,} ({p:.1f}%)', va='center', fontsize=10, weight='bold'

plt.tight_layout()
plt.show()
```

=====  
SPATIAL WORD ANALYSIS  
=====

OVERALL

Total samples:	39,129
With any spatial word:	15,445 (39.5%)
Without spatial words:	23,684 (60.5%)

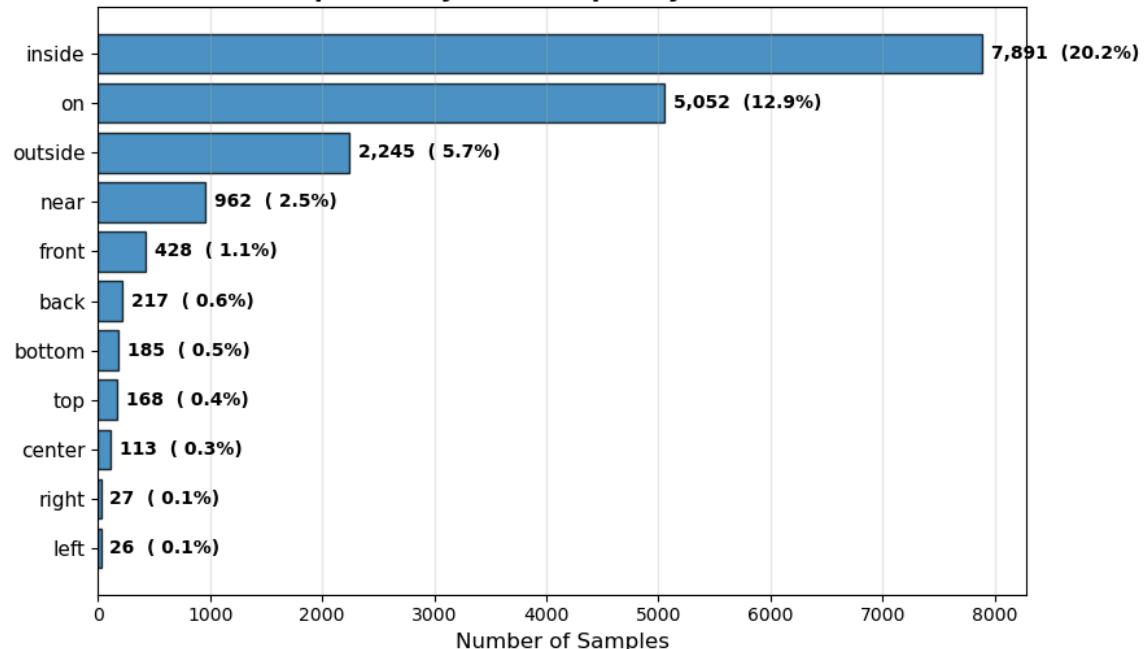
BY CLASS (any spatial word present)

Entailment :	7665/19619 ( 39.1%)
Contradiction:	7780/19510 ( 39.9%)

KEYWORD DISTRIBUTION – OVERALL

keyword	count	pct
inside	7891	20.2
on	5052	12.9
outside	2245	5.7
near	962	2.5
front	428	1.1
back	217	0.6
bottom	185	0.5
top	168	0.4
center	113	0.3
right	27	0.1
left	26	0.1

Spatial Keyword Frequency — Overall



## Observations

- Around 39.5 % of hypotheses contain at least one spatial keyword, with inside (20.2 %), on (12.9 %), and outside (5.7 %) being the most frequent.

- Because spatial language is common, models must be sensitive to positional cues in both image and text.
  - Safe augmentation: avoid geometric flips or rotations that distort spatial relations;
  - SAFE AUGMENTATIONS (preserve spatial semantics):
    - Brightness adjustment
    - Contrast adjustment
    - Gaussian noise
    - Random crop (without rotation)
    - Zoom
- RISKY AUGMENTATIONS (corrupt spatial semantics):
- geometric flips with top/bottom
  - Rotation: affects samples with top/bottom/left/right

```
In [51]: print("\n" + "="*70)
print("EDA: SEMANTIC BUCKET ANALYSIS")
print("="*70)

# Set style
sns.set_style("whitegrid")
plt.rcParams['figure.figsize'] = (14, 8)

# =====
# DEFINE ALL POTENTIAL BUCKETS
# =====

ALL_BUCKETS = {
    "gender": {"regex": r"\b(man|woman|boy|girl|male|female|gentleman|lady|m
                  "color": "#FF6B6B", "description": "Gender-specific terms"},

    "counting": {"regex": r"\b(one|two|three|four|five|six|seven|eight|\d+|cou
                  "color": "#4ECDC4", "description": "Numerical/counting terms"},

    "colors": {"regex": r"\b(blue|black|red|white|green|yellow|brown|pink|pu
                  "color": "#95E1D3", "description": "Color attributes"},

    "vehicles": {"regex": r"\b(bike|bicycle|scooter|skateboard|motorcycle|car|
                  "color": "#F38181", "description": "Vehicle types"},

    "posture": {"regex": r"\b(sit|sitting|standing|stand|talk|talking|kiss|ki
                  "color": "#AA96DA", "description": "Body postures/actions"},

    "locations": {"regex": r"\b(doorway|inside|outside|indoors|outdoors|street|
                  "color": "#FCBAD3", "description": "Location descriptors"},

    "negation": {"regex": r"\b(no|not|never|without|none|doesn't|isn't|aren't|
                  "color": "#FFFFD2", "description": "Negation words"},

    "clothing": {"regex": r"\b(shirt|shorts|pants|jeans|dress|jacket|sweater|u
                  "color": "#A8D8EA", "description": "Clothing items"},

    "activity": {"regex": r"\b(play|playing|fight|fighting|argue|arguing|inter
                  "color": "#FFAA55", "description": "Activities"},

    "spatial": {"regex": r"\b(left|right|front|behind|near|far|next to|beside|
                  "color": "#FFD3B6", "description": "Spatial relationships"},

}

print(f"\nDefined {len(ALL_BUCKETS)} semantic buckets for analysis")
```

```

# =====
# CATEGORIZE TRAINING DATA
# =====

label_col = 'label_id'

def detect_buckets(hypothesis, bucket_dict):
    """Detect which buckets a hypothesis belongs to."""
    buckets = []
    h_lower = hypothesis.lower()
    for bucket_name, bucket_info in bucket_dict.items():
        if re.search(bucket_info['regex'], h_lower):
            buckets.append(bucket_name)
    return buckets

# Categorize
df_eda = df_train.copy()
df_eda['buckets'] = df_eda['hypothesis'].apply(lambda h: detect_buckets(h, ALL_BUCKETS))
df_eda['n_buckets'] = df_eda['buckets'].apply(len)

print(f"\nProcessed {len(df_eda)} training examples")

# =====
# BUCKET FREQUENCY ANALYSIS
# =====

bucket_stats = []

for bucket_name, bucket_info in ALL_BUCKETS.items():
    df_bucket = df_eda[df_eda['buckets'].apply(lambda b: bucket_name in b)]

    if len(df_bucket) == 0:
        continue

    n_total = len(df_bucket)
    n_entail = (df_bucket[label_col] == 1).sum()
    n_contra = (df_bucket[label_col] == 0).sum()
    pct_entail = (n_entail / n_total) * 100
    pct_dataset = (n_total / len(df_eda)) * 100

    # Determine balance status
    if 45 <= pct_entail <= 55:
        balance_status = "Excellent"
        balance_color = "green"
    elif 40 <= pct_entail <= 60:
        balance_status = "Good"
        balance_color = "blue"
    elif 35 <= pct_entail <= 65:
        balance_status = "Moderate"
        balance_color = "orange"
    else:
        balance_status = "Poor"
        balance_color = "red"

    bucket_stats.append({
        'bucket': bucket_name,
        'count': n_total,
        'pct_dataset': pct_dataset,
        'entailment': n_entail,
    })

```

```

        'contradiction': n_contra,
        'pct_entailment': pct_entail,
        'balance_status': balance_status,
        'balance_color': balance_color,
        'color': bucket_info['color'],
        'description': bucket_info['description']
    })

df_bucket_stats = pd.DataFrame(bucket_stats).sort_values('count', ascending=False)

print("\nBucket Statistics (sorted by frequency):")
print("*100")
print(f"{'Bucket':<12} {'Count':<8} {'% Data':<8} {'Entail':<8} {'Contra':<8} {'")
print("-*100")

for _, row in df_bucket_stats.iterrows():
    print(f"{row['bucket']:<12} {row['count']:<8} {row['pct_dataset']:<8.1f} "
          f"{row['entailment']:<8} {row['contradiction']:<8} {row['pct_entailmen"
          f"{row['balance_status']:<12}")

# =====
# VISUALIZATION 1: BUCKET FREQUENCY
# =====

fig, axes = plt.subplots(2, 2, figsize=(16, 12))
fig.suptitle('Semantic Bucket Analysis for Oversampling Strategy', fontsize=16,

# Plot 1: Bucket frequency
ax1 = axes[0, 0]
bars1 = ax1.bars(df_bucket_stats['bucket'], df_bucket_stats['count'],
                  color=df_bucket_stats['color'].tolist())
ax1.set_xlabel('Number of Examples', fontsize=11, fontweight='bold')
ax1.set_ylabel('Semantic Bucket', fontsize=11, fontweight='bold')
ax1.set_title('(A) Bucket Frequency in Training Set', fontsize=12, fontweight='b
ax1.grid(axis='x', alpha=0.3)

# Add value labels
for i, (bucket, count) in enumerate(zip(df_bucket_stats['bucket'], df_bucket_st
    ax1.text(count + 200, i, f'{count:,}', {"df_bucket_stats.iloc[i]['pct_dataset'
        va='center', fontsize=9)

# Highlight selected buckets
selected_buckets = ['gender', 'counting']
for i, bucket in enumerate(df_bucket_stats['bucket']):
    if bucket in selected_buckets:
        bars1[i].set_edgecolor('black')
        bars1[i].set_linewidth(3)

# Plot 2: Class balance per bucket
ax2 = axes[0, 1]
x_pos = np.arange(len(df_bucket_stats))
width = 0.35

bars_entail = ax2.bar(x_pos - width/2, df_bucket_stats['entailment'], width,
                      label='Entailment', color="#4CAF50", alpha=0.8)
bars_contra = ax2.bar(x_pos + width/2, df_bucket_stats['contradiction'], width,
                      label='Contradiction', color="#F44336", alpha=0.8)

ax2.set_xlabel('Semantic Bucket', fontsize=11, fontweight='bold')
ax2.set_ylabel('Number of Examples', fontsize=11, fontweight='bold')
ax2.set_title('(B) Class Distribution per Bucket', fontsize=12, fontweight='bold'

```

```

ax2.set_xticks(x_pos)
ax2.set_xticklabels(df_bucket_stats['bucket'], rotation=45, ha='right')
ax2.legend(fontsize=10)
ax2.grid(axis='y', alpha=0.3)

# Highlight selected buckets
for i, bucket in enumerate(df_bucket_stats['bucket']):
    if bucket in selected_buckets:
        ax2.axvspan(i-0.4, i+0.4, alpha=0.1, color='gold', zorder=0)

# Plot 3: Entailment percentage
ax3 = axes[1, 0]
colors_balance = df_bucket_stats['balance_color'].tolist()
bars3 = ax3.barh(df_bucket_stats['bucket'], df_bucket_stats['pct_entailment'],
                  color=colors_balance, alpha=0.7)

# Add reference lines
ax3.axvline(50, color='green', linestyle='--', linewidth=2, alpha=0.5, label='P')
ax3.axvline(45, color='orange', linestyle='--', linewidth=1, alpha=0.3)
ax3.axvline(55, color='orange', linestyle='--', linewidth=1, alpha=0.3)

ax3.set_xlabel('Entailment Percentage (%)', fontsize=11, fontweight='bold')
ax3.set_ylabel('Semantic Bucket', fontsize=11, fontweight='bold')
ax3.set_title('(C) Class Balance: Entailment %', fontsize=12, fontweight='bold')
ax3.set_xlim(0, 100)
ax3.legend(fontsize=9)
ax3.grid(axis='x', alpha=0.3)

# Add value labels
for i, (bucket, pct) in enumerate(zip(df_bucket_stats['bucket'], df_bucket_stats['pct_entailment'])):
    status = df_bucket_stats.iloc[i]['balance_status']
    ax3.text(pct + 2, i, f'{pct:.1f}% ({status})', va='center', fontsize=8)

# Highlight selected buckets
for i, bucket in enumerate(df_bucket_stats['bucket']):
    if bucket in selected_buckets:
        bars3[i].set_edgecolor('black')
        bars3[i].set_linewidth(3)

# Plot 4: Selection criteria scatter
ax4 = axes[1, 1]

# Create scatter
scatter = ax4.scatter(df_bucket_stats['count'], df_bucket_stats['pct_entailment'],
                      s=df_bucket_stats['pct_dataset']*50,
                      c=range(len(df_bucket_stats)),
                      cmap='tab10',
                      alpha=0.6,
                      edgecolors='black',
                      linewidth=1)

# Highlight selected buckets
for _, row in df_bucket_stats.iterrows():
    if row['bucket'] in selected_buckets:
        ax4.scatter(row['count'], row['pct_entailment'],
                    s=row['pct_dataset']*50,
                    edgecolors='gold',
                    linewidth=4,
                    facecolors='none',
                    marker='o',
                    alpha=0.6)

```

```

zorder=10)

# Add Labels
for _, row in df_bucket_stats.iterrows():
    offset = 1.5 if row['bucket'] in selected_buckets else 1.0
    ax4.annotate(row['bucket'],
                 (row['count'], row['pct_entailment']),
                 xytext=(5, 5*offset),
                 textcoords='offset points',
                 fontsize=9,
                 fontweight='bold' if row['bucket'] in selected_buckets else 'normal',
                 bbox=dict(boxstyle='round', pad=0.3, facecolor='yellow', alpha=0.5))

# Reference zones
ax4.axhspan(40, 60, alpha=0.1, color='green', label='Balanced Zone (40-60%)')
ax4.axhline(50, color='green', linestyle='--', linewidth=1, alpha=0.5)

ax4.set_xlabel('Number of Examples', fontsize=11, fontweight='bold')
ax4.set_ylabel('Entailment Percentage (%)', fontsize=11, fontweight='bold')
ax4.set_title('(D) Selection Criteria: Size vs Balance', fontsize=12, fontweight='bold')
ax4.legend(fontsize=9, loc='upper right')
ax4.grid(alpha=0.3)

plt.tight_layout()
plt.savefig(RUNS_DIR / 'bucket_analysis_oversampling.png', dpi=300, bbox_inches='tight')
print(f"\nSaved: {RUNS_DIR / 'bucket_analysis_oversampling.png'}")
plt.show()

```

=====
EDA: SEMANTIC BUCKET ANALYSIS
=====

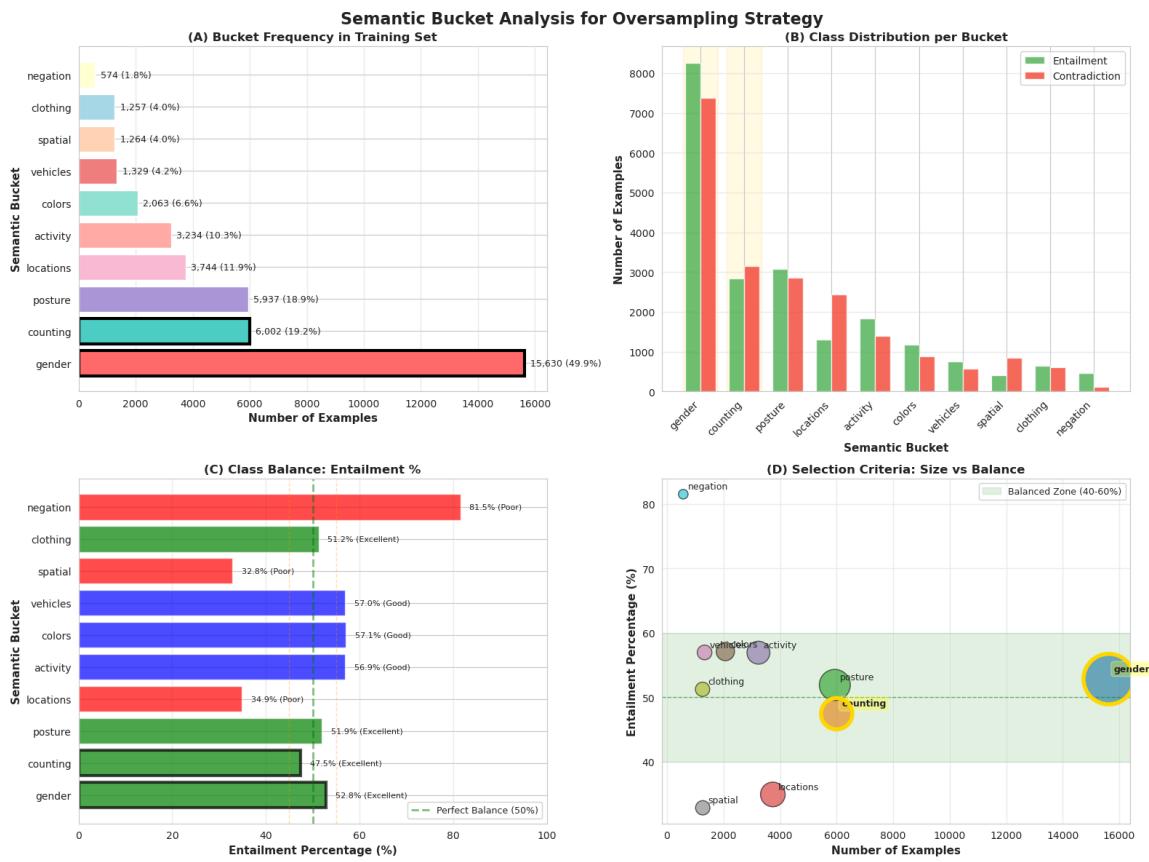
Defined 10 semantic buckets for analysis

Processed 31,340 training examples

Bucket Statistics (sorted by frequency):

Bucket	Count	% Data	Entail	Contra	% Entail	Balance
<hr/>						
gender	15630	49.9	8256	7374	52.8	Excellent
counting	6002	19.2	2848	3154	47.5	Excellent
posture	5937	18.9	3082	2855	51.9	Excellent
locations	3744	11.9	1307	2437	34.9	Poor
activity	3234	10.3	1841	1393	56.9	Good
colors	2063	6.6	1178	885	57.1	Good
vehicles	1329	4.2	757	572	57.0	Good
spatial	1264	4.0	415	849	32.8	Poor
clothing	1257	4.0	644	613	51.2	Excellent
negation	574	1.8	468	106	81.5	Poor

Saved: /home/ec2-user/SageMaker/abeyt/runs/bucket\_analysis\_oversampling.png



Based on commonly occurring words, we did an analysis of hypothesis distribution. We could see words related to negation, gender, colours, counting etc are very common. Some are balanced and some are unbalanced in the data. So we will do a detailed analysis on the further modelling process to tackle this problem.

## Auxiliary Dataset: SNLI-VE (proxy warmup)

100,000 balanced samples from Flickr30k images Used for transfer learning (Phase-A warmup)

## FLICKR30k Dataset

Our A2 dataset, while reasonably sized at approximately 38,000 image-hypothesis pairs with balanced 50-50 class distribution (entailment and contradiction), presents a challenge for training deep neural networks from scratch given the massive parameter count of modern vision-language models. To provide better initialization than random weights, we leverage transfer learning from a related task using another dataset.

The Flickr30k dataset is a widely-used benchmark in vision-language research, containing 31,783 images collected from Flickr with each image accompanied by five human-written captions describing the visual content—totaling over 150,000 image-caption pairs that capture diverse real-world scenes, objects, and human activities. Building upon Flickr30k, the SNLI-VE (Stanford Natural Language Inference - Visual Entailment) dataset was

created by combining Flickr30k images with hypothesis sentences from the Stanford Natural Language Inference corpus, resulting in a large-scale visual entailment benchmark with approximately 565,000 image-hypothesis pairs across three labels: entailment, contradiction, and neutral.

We select SNLI-VE as our proxy warmup task because it directly addresses the same core problem—determining whether a textual statement is supported or contradicted by visual evidence—making it the closest available large-scale dataset to our target A2 task. After filtering to match our binary classification setup (removing neutral samples and balancing classes), we utilize a 100,000-sample subset of SNLI-VE to pretrain our task-specific head (cross-attention, fusion, and classifier layers), which learns general visual-linguistic reasoning patterns before specializing to the A2 domain. This transfer learning strategy provides our model with exposure to diverse visual scenes and reasoning patterns from Flickr30k's rich imagery, significantly improving initialization compared to random weights and reducing the risk of overfitting on our smaller target dataset.

let's have a look on Flicker30K dataset

```
In [49]: #Flicker30k EDA

# Point BASE to where 'abeyt' folder lives
BASE = Path("/home/ec2-user/SageMaker/abeyt") # e.g., Path("/home/ec2-user/Sage

IMG_DIR    = BASE / "Flickr"
JSON_PATH  = BASE / "dataset_flickr30k.json"

assert IMG_DIR.exists(), f"Image directory not found: {IMG_DIR.resolve()}"
assert JSON_PATH.exists(), f"JSON file not found: {JSON_PATH.resolve()}""

# Load captions: {"1000092795.jpg": ["<start> ... <end>", ... x5]}
with open(JSON_PATH, "r", encoding="utf-8") as f:
    captions = json.load(f)

# Find images present on disk AND in JSON
image_files = sorted([p for p in IMG_DIR.iterdir() if p.suffix.lower() in (".jpg")]
pairs = [(p, captions[p.name]) for p in image_files if p.name in captions]

# Take first 4 for display
pairs = pairs[:4]

# 1) Grid: image + title as the image id only (no caption)
fig, axes = plt.subplots(2, 2, figsize=(12, 9))
axes = axes.flatten()
for ax, (img_path, caps) in itertools.zip_longest(axes, pairs, fillvalue=None):
    ax.axis("off")
    if img_path is None:
        continue
    img = Image.open(img_path).convert("RGB")
    ax.imshow(img)
    ax.set_title(img_path.name, fontsize=11) # title = image id only

plt.tight_layout()
plt.show()

# 2) Table: all 5 captions (cleaned) for those images
```

```
def clean(c):
    return c.replace("<start>", "").replace("<end>", "").strip()

rows = []
for img_path, caps in pairs:
    cap_clean = [clean(c) for c in caps]
    rows.append({
        "image": img_path.name,
        "caption_1": cap_clean[0],
        "caption_2": cap_clean[1],
        "caption_3": cap_clean[2],
        "caption_4": cap_clean[3],
        "caption_5": cap_clean[4],
    })

pd.DataFrame(rows)
```



134206.jpg



36979.jpg



65567.jpg



81641.jpg

Out[49]:

	image	caption_1	caption_2	caption_3	caption_4	caption_5
0	134206.jpg	The players of the baseball team are standing ...	Baseball players are playing on a field in a s...	A team plays baseball at a large crowded stadium.	A crowd cheers on a baseball team.	Game is playing in the stadium
1	36979.jpg	A group of friends playing cards and trying to...	A group of college students gathers to play te...	Several men play cards while around a green ta...	A group of several men playing poker.	Six white males playing poker.
2	65567.jpg	A bearded man, and a girl in a red dress are g...	The group of people are assembling for a wedding.	A man and woman dressed for a wedding function.	A woman holds a man's arm at a formal event.	A wedding party walks out of a building.
3	81641.jpg	A man does acrobatics outside of a middle east...	A man is upside down on an outside gymnastics ...	The man is flipping over the bar.	A man is upside-down on a pole.	A man doing gymnastic stunts

Now Lets check for the SNLI VE dataset

In [52]:

```
# ===== SNLI-VE quick EDA: show images with [label] hypothesis =====
# --- Paths (adjust if s differ) ---
DATA_DIR      = "/home/ec2-user/SageMaker/abeyt"
BINARY_PATH   = Path(DATA_DIR) / "runs/snli_ve_warmup_keras_hub/snli_ve_binary"

# --- Load filtered binary splits (0: entailment, 1: contradiction) ---
def drop_neutral(ex): return ex["label"] != 1
def remap_binary(ex): ex["label"] = 0 if ex["label"] == 0 else 1; return ex

try:
    dd = load_from_disk(str(BINARY_PATH))
    train_raw, val_raw, test_raw = dd["train"], dd["validation"], dd["test"]
    print("Loaded SNLI-VE (binary) from:", BINARY_PATH)
except Exception:
    print("Filtered splits not found on disk. Recreating from HuggingFace (one-t")
    ds = load_dataset("HuggingFaceM4/SNLI-VE")
    train_raw = ds["train"].filter(drop_neutral).map(remap_binary).flatten_indices
    val_raw   = ds["validation"].filter(drop_neutral).map(remap_binary).flatten_indices
    test_raw  = ds["test"].filter(drop_neutral).map(remap_binary).flatten_indices

ID2LABEL = {0: "entailment", 1: "contradiction"}

def show_samples(split, n=8, cols=4, title="SNLI-VE samples"):
    n = min(n, len(split))
    idx = np.random.default_rng(42).choice(len(split), size=n, replace=False)
    rows = (n + cols - 1) // cols

    fig, axes = plt.subplots(rows, cols, figsize=(5*cols, 4.5*rows))
    axes = np.atleast_1d(axes).ravel()

    for ax, i in itertools.zip_longest(axes, idx, fillvalue=None):
        ax.axis("off")
        if i is None:
            continue
        ex = split[int(i)]
```

```

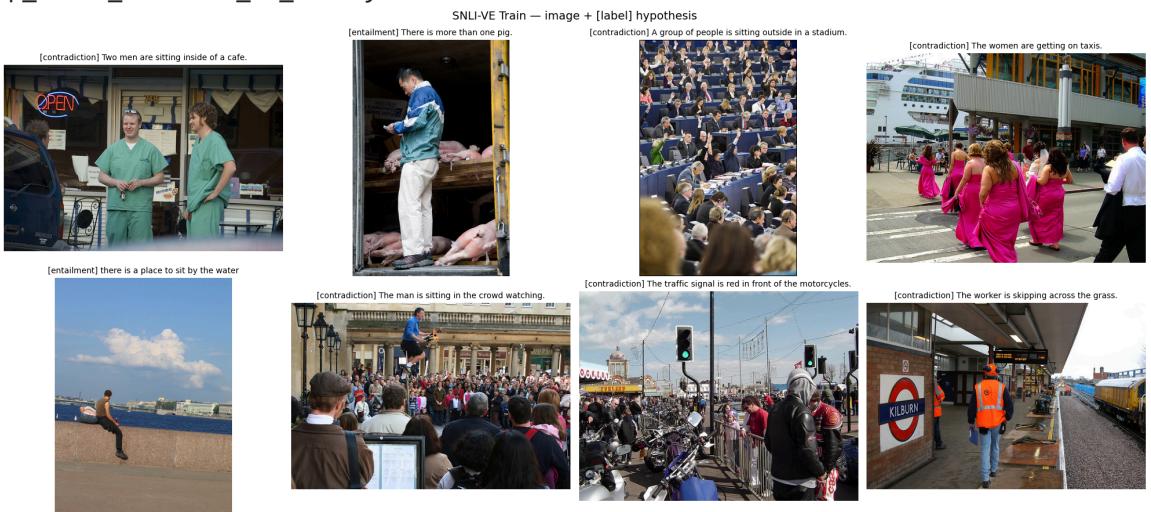
# 'image' is already a PIL.Image in this dataset
img: Image.Image = ex["image"].convert("RGB")
ax.imshow(img)
hyp = str(ex["hypothesis"])
lab = ID2LABEL[int(ex["label"])]
ax.set_title(f"[{lab}] " + textwrap.fill(hyp, width=60), fontsize=10)

fig.suptitle(title, fontsize=14)
plt.tight_layout()
plt.show()

# Examples (pick the split you want to preview)
show_samples(train_raw, n=8, cols=4, title="SNLI-VE Train - image + [label] hypo")
# show_samples(val_raw, n=8, cols=4, title="SNLI-VE Val - image + [Label] hypo")
# show_samples(test_raw, n=8, cols=4, title="SNLI-VE Test - image + [label] hypo")

```

Loading dataset from disk: 0% | 0/101 [00:00<?, ?it/s]  
 Loaded SNLI-VE (binary) from: /home/ec2-user/SageMaker/abeyt/runs/snli\_ve\_warmup\_keras\_hub/snli\_ve\_binary



## Target Setting

Achieve **80% Performance with Strong Generalization** Our primary objective is to develop a visual entailment classifier that achieves  $\geq 80\%$  on both **accuracy and Macro-F1** score on held-out test data, with strong evidence of generalization to unseen distributions. We intentionally target both metrics equally(as bdataet is balanced): accuracy provides an intuitive overall performance measure, while Macro-F1 ensures balanced performance across entailment and contradiction classes Practical Utility: 80% accuracy makes the system viable for human-in-the-loop applications where model predictions can be reviewed.

Rather than pursuing a single end-to-end training run, we adopt an incremental improvement methodology that systematically addresses different aspects of model performance: The staged approach provides interpretability (each phase's contribution is measurable), risk mitigation (failures in later phases don't invalidate earlier work), and computational efficiency (early stopping if targets are met sooner than expected).

## Model Architecture

## SNLI-VE Warm-up Architecture Summary

Our **SNLI-VE warm-up model** employs a simplified architecture designed for efficient transfer learning with **frozen pretrained backbones**.

- **Vision Encoder:** ViT-Base/16 pretrained on ImageNet-21k encodes each image into **196 patch embeddings** ( $14 \times 14$  grid).
- **Text Encoder:** DistilBERT processes hypothesis sentences into **24 token embeddings**.
- Both encoders remain **frozen** to preserve pretrained knowledge and prevent overfitting.

### Fusion Mechanism:

- The `[CLS]` token from DistilBERT's output serves as a **compact text representation**.
- It is used as the **query** in a **12-head cross-attention module** that attends to all **196 visual patch tokens** from ViT.
- This allows the model to focus on visually relevant regions for verifying the hypothesis.
- The attended features pass through a **residual connection** and **Layer Normalization** for stable training.
- The fused representation is then **flattened** and passed into a classification head:
  - `Dense(512, activation='gelu')`
  - `Dropout(0.1)`
  - `Dense(1) → binary logits output`

### Training Summary:

- Total trainable parameters: **≈2.76M** (cross-attention + classification head)
- Frozen pretrained parameters: **≈152M**
- Dataset: **100,000 SNLI-VE balanced samples** (50k entailment + 50k contradiction)
- Loss: **BinaryCrossentropy(from\_logits=True)**
- Optimizer: **AdamW(lr=2e-4, weight\_decay=1e-4)**
- Metric: **BinaryAccuracy(threshold=0.0)**
- The model achieves efficient warm-up and strong multimodal alignment despite limited trainable capacity.

## Component Rationale

### Why ViT-Base (Vision Transformer)?

**Choice:** `vit_base_patch16_224_imagenet21k`

**Rationale:**

- **State-of-the-Art Performance:**  
ViT has become the standard for vision tasks, outperforming CNNs on many benchmarks.
- **Attention-Based:**  
Naturally suited for cross-modal reasoning — attention focuses on relevant image regions.
- **Patch-Based Representation:**
  - Divides a  $224 \times 224$  image into **14x14 = 196 patches** (16x16 pixels each).
  - Each patch becomes a “visual word” the text can attend to.
  - (CLS token exists internally but not explicitly used here.)
- **ImageNet-21k Pretraining:**  
Trained on 21,000 categories → rich, diverse visual representations.
- **Frozen Backbone:**
  - $\approx 86M$  parameters — too expensive to fine-tune.
  - Pretrained features already capture strong object and scene semantics.
  - Freezing prevents catastrophic forgetting.

### Why NOT CNN (ResNet/EfficientNet)?

- CNNs output spatial feature maps that are harder to align with sequential text.
  - ViT’s **token-based output** naturally pairs with BERT’s **token-based embeddings** for seamless fusion.
- 

### Why DistilBERT (Text Encoder)?

**Choice:** `distil_bert_base_en_uncased`

#### Rationale:

- **Efficiency:**  
40% smaller, 60% faster than BERT-Base while retaining ~97% of its performance.
- **Strong Language Understanding:**
  - Pretrained on **BookCorpus + Wikipedia (3.3B words)**.
  - Understands syntax, semantics, and contextual nuance.
- **Handles Complex Reasoning:**  
Supports logical negations (“not”, “without”) and spatial relations (“left of”, “behind”).
- **Sequence Output:**  
Returns `[BATCH, SEQ_LEN=24, 768]` hidden states → fine-grained token-level representations.  
The `[CLS]` token is extracted as the text summary vector.
- **Frozen Backbone:**
  - $\approx 66M$  parameters.
  - Pretrained knowledge sufficient for this task; freezing improves stability and efficiency.

### Why NOT Smaller Models (DistilRoBERTa, TinyBERT)?

- Hypotheses are short (8–10 words) but semantically complex.
- Require robust reasoning for negation and spatial language.

### Why NOT Larger Models (BERT-Large, RoBERTa)?

- High computational cost with marginal accuracy gain.
  - DistilBERT offers the best **accuracy–efficiency trade-off**.
- 

## Why Cross-Attention (Text → Image)?

**Design:** Multi-Head Attention where **text queries image features**.

### Rationale:

- **Alignment Mechanism:**
  - Text tokens (specifically `[CLS]`) attend to relevant image patches.
  - Example: “A dog on the left” → attention focuses on left-side patches with the dog.
- **Flexible Reasoning:**
  - **12 attention heads**, each with **key\_dim = 64**, and **dropout = 0.1**.
  - Different heads specialize in distinct relationships:
    - Head 1: Object presence (“is there a dog?”)
    - Head 2: Spatial relation (“where is it?”)
    - Head 3+: Contextual alignment
- **Direction Matters:**
  - **Text → Image**: “Given this hypothesis, what regions in the image support or contradict it?”
  - Works better than bidirectional fusion for entailment-style reasoning.
- **Trainable with Frozen Backbones:**
  - Only cross-attention and classification head (~2.7M params) are trainable.
  - Allows efficient task-specific alignment without updating heavy backbones.

### Architecture Details (from code):

- Attention heads: **12**
- Key dimension: **64**
- Dropout: **0.1**
- Residual connection: `x = cls_q + attn`
- Layer normalization: `epsilon = 1e-6`
- Classification head: `Dense(512, 'gelu') → Dropout(0.1) → Dense(1)`

### Note:

There is **no gated fusion** (no sigmoid-based gating or weighted combination). Fusion relies purely on **cross-attention + residual addition + normalization**, not on learned gating.

---

## Additional Notes

- With Assignment-2's **~38k samples**, we still leverage ViT-B/16 and DistilBERT for strong pretrained priors, improving **out-of-distribution generalization**.
- The **text→image cross-attention** enables hypotheses to attend to salient visual regions, reducing the risk of memorizing dataset biases.
- For **OOV (Out-of-Vocabulary)** words, **WordPiece subword tokenization** decomposes rare/unseen words into known subunits, ensuring meaningful embeddings and maintaining performance even for novel entities or morphology.

Given below is a summary of the SNLI-VE Warm-up Architecture and training setup

```
In [83]: # =====
# PHASE 0: SNLI-VE Warm-up Architecture & Results
# =====

print("""
        SNLI-VE WARM-UP ARCHITECTURE

1. INPUT LAYER

    • Image: (224, 224, 3) – RGB image
    • Text: (SEQ_LEN=24,) – tokenized hypothesis
        └ token_ids: [CLS] + tokens + [SEP] + [PAD]
        └ padding_mask: attention mask

2. VISION ENCODER (FROZEN )

    • Backbone: ViT-Base/16 (vit_base_patch16_224_imagenet21k)
    • Output: Patch embeddings → (batch, 196, 768)
    • Status: FROZEN (no weight updates in warm-up)

3. TEXT ENCODER (FROZEN )

    • Backbone: DistilBERT (distil_bert_base_en_uncased)
    • Pretrained: English Wikipedia + BookCorpus
    • Output: Token embeddings → (batch, 24, 768)
        └ sequence_output: all token embeddings
        └ [CLS] token: position 0
    • Status: FROZEN (no weight updates in warm-up)

4. FUSION MODULE (TRAINABLE )

    a) Extract [CLS] token from text: (batch, 1, 768)

    b) Cross-Attention (Text queries Image):
        • MultiHeadAttention:
            └ num_heads: 12
            └ key_dim: 64
            └ dropout: 0.1
            └ Query: [CLS] token (text)
            └ Key/Value: patch embeddings (vision)
            └ Output: attended visual features → (batch, 1, 768)

    c) Residual Connection + Layer Normalization:
        • x = [CLS] + attention_output
```

```
• x = LayerNorm(x, epsilon=1e-6)
```

## 5. CLASSIFICATION HEAD (TRAINABLE )

- Flatten: (batch, 1, 768) → (batch, 768)
- Dense(512, activation='gelu')
- Dropout(0.1)
- Dense(1) – logits (binary classification)

## 6. TRAINING CONFIGURATION

- Dataset: SNLI-VE 100k balanced (50k entail + 50k contradict)
- Loss: BinaryCrossentropy(from\_logits=True)
- Optimizer: AdamW(lr=2e-4, weight\_decay=1e-4, clipnorm=1.0)
- Metrics: BinaryAccuracy(threshold=0.0)
- Epochs: 6 (with early stopping, patience=1)
- Batch size: 32

## OUTPUT

- Shape: (batch, 1) – raw logits
- Saved artifact: best.weights.h5 (2.76M params)

```
""")
```

## SNLI-VE WARM-UP ARCHITECTURE

## 1. INPUT LAYER

- Image: (224, 224, 3) – RGB image
- Text: (SEQ\_LEN=24,) – tokenized hypothesis
  - | token\_ids: [CLS] + tokens + [SEP] + [PAD]
  - | padding\_mask: attention mask

## 2. VISION ENCODER (FROZEN )

- Backbone: ViT-Base/16 (vit\_base\_patch16\_224\_imagenet21k)
- Output: Patch embeddings → (batch, 196, 768)
- Status: FROZEN (no weight updates in warm-up)

## 3. TEXT ENCODER (FROZEN )

- Backbone: DistilBERT (distil\_bert\_base\_en\_uncased)
- Pretrained: English Wikipedia + BookCorpus
- Output: Token embeddings → (batch, 24, 768)
  - | sequence\_output: all token embeddings
  - | [CLS] token: position 0
- Status: FROZEN (no weight updates in warm-up)

## 4. FUSION MODULE (TRAINABLE )

- a) Extract [CLS] token from text: (batch, 1, 768)
- b) Cross-Attention (Text queries Image):
  - MultiHeadAttention:
    - | num\_heads: 12
    - | key\_dim: 64
    - | dropout: 0.1
    - | Query: [CLS] token (text)
    - | Key/Value: patch embeddings (vision)
    - | Output: attended visual features → (batch, 1, 768)
- c) Residual Connection + Layer Normalization:
  - $x = [\text{CLS}] + \text{attention\_output}$
  - $x = \text{LayerNorm}(x, \text{epsilon}=1e-6)$

## 5. CLASSIFICATION HEAD (TRAINABLE )

- Flatten: (batch, 1, 768) → (batch, 768)
- Dense(512, activation='gelu')
- Dropout(0.1)
- Dense(1) – logits (binary classification)

## 6. TRAINING CONFIGURATION

- Dataset: SNLI-VE 100k balanced (50k entail + 50k contradict)
- Loss: BinaryCrossentropy(from\_logits=True)
- Optimizer: AdamW(lr=2e-4, weight\_decay=1e-4, clipnorm=1.0)
- Metrics: BinaryAccuracy(threshold=0.0)
- Epochs: 6 (with early stopping, patience=1)
- Batch size: 32

## OUTPUT

- Shape: (batch, 1) – raw logits
- Saved artifact: best.weights.h5 (2.76M params)

The classification head takes the CLS token after cross-attention, flattens it, passes it through a 512-unit GELU layer with dropout, and finally outputs a single scalar logit through a Dense(1) layer with no activation. We use logits instead of an explicit sigmoid because BinaryCrossentropy(from\_logits=True) internally applies a numerically stable sigmoid, gives better gradients, and allows flexible thresholding ( $\tau$ -sweep) later. During evaluation, we apply a sigmoid and tune the decision threshold, while for metrics a threshold of 0.0 on logits corresponds to 0.5 probability.

We will load the warm up model in this notebook. The warm up training of the dataset is attached. see appendix Let's load the model and see the performance and learning curves

```
In [17]: # Load and display history
with open(SNLI_HIST, 'r') as f:
    snli_history = json.load(f)

n_epochs = len(snli_history.get('loss', []))
print(f"\nEpochs trained: {n_epochs}")
print(f"\n{'Epoch':<8} {'Train Loss':<12} {'Val Loss':<12} {'Train Acc':<12} {'Val Acc':<12}\n")
print("-"*60)

for i in range(n_epochs):
    print(f"{i+1:<8} "
          f"{snli_history['loss'][i]:<12.4f} "
          f"{snli_history['val_loss'][i]:<12.4f} "
          f"{snli_history['acc'][i]:<12.4f} "
          f"{snli_history['val_acc'][i]:<12.4f}")

print("\n Best epoch: {np.argmin(snli_history['val_loss']) + 1} "
      f"(val_loss={min(snli_history['val_loss']):.4f})")
print(f" Final validation accuracy: {snli_history['val_acc'][-1]:.4f}\n")

print("\n" + "-"*70)
print(" LEARNING CURVES")
print("-"*70)

# Plot curves
plot_history_curves(SNLI_HIST, title="SNLI-VE Warmup (100k Balanced)", out_dir=SNLI_WTS)

print(f" Warm-up weights : {SNLI_WTS}")
```

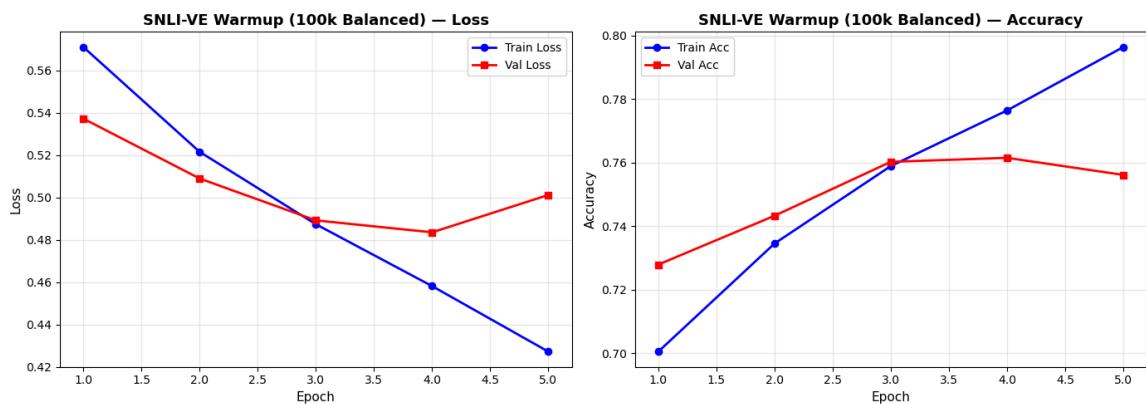
Epochs trained: 5

Epoch	Train Loss	Val Loss	Train Acc	Val Acc
1	0.5711	0.5373	0.7005	0.7279
2	0.5216	0.5090	0.7345	0.7433
3	0.4874	0.4892	0.7589	0.7602
4	0.4582	0.4836	0.7764	0.7615
5	0.4272	0.5013	0.7964	0.7561

Best epoch: 4 (val\_loss=0.4836)

Final validation accuracy: 0.7561

#### LEARNING CURVES



Warm-up weights : /home/ec2-user/SageMaker/abeyt/runs/snli\_ve\_warmup\_keras\_hu/b/sub100k\_balanced\_e6/best.weights.h5

The model achieved a best validation accuracy of 76.15% and validation loss of 0.4836 at Epoch 4, with stable training curves and minimal overfitting. This establishes a strong baseline for subsequent transfer learning and fine-tuning on A2 data.

## Phase 1: A2 Stage-0 (Frozen Backbones)

Warm-start A2 using the SNLI-VE-trained fusion head while keeping both encoders frozen. This adapts only the small classification head to A2, establishes a baseline, and avoids overfitting before We kept effectively the same architecture.No structural change: same inputs, same ViT/DistilBERT backbones, same cross-attention (CLS→ViT tokens) + residual + LayerNorm, and same head ('Flatten → Dense(512, GELU) → Dropout → Dense(1) logit).

We performed grouped splitting by image\_id using a two-stage GroupShuffleSplit (80/20 train+val/test, then 80/20 train/val within train+val) while keeping all captions of an image in the same split. This prevents leakage (no shared images across splits) and provides a more realistic evaluation of generalization to unseen images, which we verified by checking zero overlap in image\_id sets.

In [7]: # 1A: Create/Load Grouped Splits

```

if SPLIT_JSON.exists():
    print(f" Loading existing splits from: {SPLIT_JSON}")
    with open(SPLIT_JSON, 'r') as f:
        splits = json.load(f)

    train_idx = np.array(splits["train_idx"])
    val_idx = np.array(splits["val_idx"])
    test_idx = np.array(splits["test_idx"])

else:
    print("Creating new grouped splits...")

    # First split: train+val vs test (80/20)
    gss_test = GroupShuffleSplit(n_splits=1, test_size=0.20, random_state=SEED)
    trainval_idx, test_idx = next(gss_test.split(
        X=df.index,
        groups=df["image_id"]
    ))

    # Second split: train vs val from trainval (80/20 of trainval)
    df_trainval = df.iloc[trainval_idx]
    gss_val = GroupShuffleSplit(n_splits=1, test_size=0.20, random_state=SEED)
    train_sub_idx, val_sub_idx = next(gss_val.split(
        X=df_trainval.index,
        groups=df_trainval["image_id"]
    ))

    train_idx = trainval_idx[train_sub_idx]
    val_idx = trainval_idx[val_sub_idx]

    # Save splits
    splits = {
        "train_idx": train_idx.tolist(),
        "val_idx": val_idx.tolist(),
        "test_idx": test_idx.tolist(),
        "seed": SEED,
        "created_at": time.strftime("%Y-%m-%d %H:%M:%S")
    }

    with open(SPLIT_JSON, 'w') as f:
        json.dump(splits, f, indent=2)

# Create split DataFrames
df_train = df.iloc[train_idx].copy().reset_index(drop=True)
df_val = df.iloc[val_idx].copy().reset_index(drop=True)
df_test = df.iloc[test_idx].copy().reset_index(drop=True)

print(f"\nSplit sizes:")
print(f" Train: {len(df_train)} samples, {df_train['image_id'].nunique()}")
print(f" Val: {len(df_val)} samples, {df_val['image_id'].nunique()}")
print(f" Test: {len(df_test)} samples, {df_test['image_id'].nunique()}")


# Verify no image Leakage
train_imgs = set(df_train["image_id"])
val_imgs = set(df_val["image_id"])
test_imgs = set(df_test["image_id"])

print(" No image leakage detected!")

```

```
Loading existing splits from: /home/ec2-user/SageMaker/abeyt/runs/_shared/a2_grouped_split.json
```

Split sizes:

```
Train: 31,340 samples, 15,658 unique images
Val: 3,922 samples, 1,957 unique images
Test: 3,867 samples, 1,958 unique images
No image leakage detected!
```

We tokenize hypotheses with DistilBERT's WordPiece preprocessor (sequence\_length = 24), producing input\_ids and attention\_mask; subword tokenization robustly handles OOV terms by splitting them into known pieces. To keep runs reproducible and fast, we cache tokenized outputs per split in compressed NPZ files (train/val/test), validating keys on load and regenerating any invalid caches.

In [8]:

```
# =====
# 1B: Tokenize Text (with caching)
# =====

TOK_CACHE_DIR = RUNS_DIR / "a2_tokens_cache"
TOK_CACHE_DIR.mkdir(exist_ok=True)

TOK_FILES = {
    "train": TOK_CACHE_DIR / "train_tokens.npz",
    "val": TOK_CACHE_DIR / "val_tokens.npz",
    "test": TOK_CACHE_DIR / "test_tokens.npz",
}

# Initialize tokenizer
Tokenizer = keras_hub.models.DistilBertTokenizer
Preprocessor = keras_hub.models.DistilBertPreprocessor

tokenizer = Tokenizer.from_preset(DISTIL_PRESET)
preprocessor = Preprocessor.from_preset(DISTIL_PRESET, sequence_length=SEQ_LEN)

def tokenize_split_to_npz(split_df, out_path: Path, bs=4096):
    """Tokenize text and save to NPZ cache."""
    if out_path.exists():
        # Verify cache has correct keys
        try:
            z = np.load(out_path)
            if "input_ids" in z and "attention_mask" in z and "labels" in z:
                print(f" Using cached: {out_path.name}")
                z.close()
                return
            else:
                print(f" Invalid cache, regenerating: {out_path.name}")
                z.close()
                out_path.unlink() # Delete corrupt cache
        except Exception as e:
            print(f" Cache read error ({e}), regenerating: {out_path.name}")
            try:
                out_path.unlink()
            except:
                pass

    print(f" Tokenizing {len(split_df)} samples...")
    texts = split_df["hypothesis"].tolist()
```

```

labels = split_df["label_id"].values.astype(np.int32)

# Batch tokenize
ds = (tf.data.Dataset.from_tensor_slices(texts)
      .batch(bs)
      .map(lambda x: preprocess(x), num_parallel_calls=tf.data.AUTOTUNE)
      .prefetch(tf.data.AUTOTUNE))

ids_chunks, mask_chunks = [], []
for out in ds:
    ids_chunks.append(out["token_ids"])
    mask_chunks.append(out["padding_mask"])

ids = tf.concat(ids_chunks, axis=0).numpy().astype(np.int32)
mask = tf.concat(mask_chunks, axis=0).numpy().astype(np.int32)

np.savez_compressed(out_path, input_ids=ids, attention_mask=mask, labels=labels)
print(f" Saved: {out_path.name}")

# Tokenize all splits
print("\nProcessing splits:")
for split_name, split_df in [("train", df_train), ("val", df_val), ("test", df_test)]:
    tokenize_split_to_npz(split_df, TOK_FILES[split_name])

# Load tokenized data
def load_tokens(path: Path):
    """Load tokenized data from NPZ cache."""
    z = np.load(path)
    return z["input_ids"], z["attention_mask"], z["labels"]

ids_train, mask_train, y_train = load_tokens(TOK_FILES["train"])
ids_val, mask_val, y_val = load_tokens(TOK_FILES["val"])
ids_test, mask_test, y_test = load_tokens(TOK_FILES["test"])

print(f"\nToken shapes:")
print(f" Train: {ids_train.shape}")
print(f" Val: {ids_val.shape}")
print(f" Test: {ids_test.shape}")

```

Downloading from [https://www.kaggle.com/api/v1/models/keras/distil\\_bert/keras/distil\\_bert\\_base\\_en\\_uncased/3/download/config.json...](https://www.kaggle.com/api/v1/models/keras/distil_bert/keras/distil_bert_base_en_uncased/3/download/config.json...)

100%|██████████| 462/462 [00:00<00:00, 1.14MB/s]

Downloading from [https://www.kaggle.com/api/v1/models/keras/distil\\_bert/keras/distil\\_bert\\_base\\_en\\_uncased/3/download/tokenizer.json...](https://www.kaggle.com/api/v1/models/keras/distil_bert/keras/distil_bert_base_en_uncased/3/download/tokenizer.json...)

100%|██████████| 794/794 [00:00<00:00, 1.59MB/s]

Downloading from [https://www.kaggle.com/api/v1/models/keras/distil\\_bert/keras/distil\\_bert\\_base\\_en\\_uncased/3/download/assets/tokenizer/vocabulary.txt...](https://www.kaggle.com/api/v1/models/keras/distil_bert/keras/distil_bert_base_en_uncased/3/download/assets/tokenizer/vocabulary.txt...)

100%|██████████| 226k/226k [00:01<00:00, 197kB/s]

I0000 00:00:1761204194.833362 7814 gpu\_device.cc:2019] Created device /job:localhost/replica:0/task:0/device:GPU:0 with 13760 MB memory: -> device: 0, name: Tesla T4, pci bus id: 0000:00:1e.0, compute capability: 7.5

Processing splits:  
 Using cached: train\_tokens.npz  
 Using cached: val\_tokens.npz  
 Using cached: test\_tokens.npz

Token shapes:  
 Train: (31340, 24)  
 Val: (3922, 24)  
 Test: (3867, 24)

In [9]:

```
# =====
# 1C: Create tf.data Pipelines (NO augmentation in Stage-0)
# =====

def pil_to_float32(image_pil, size=IMG_SIZE):
    """Convert PIL image to normalized float32 array."""
    img = image_pil.convert("RGB").resize((size, size))
    return np.asarray(img, dtype=np.float32) / 255.0

# Signature for tf.data
sig_input = {
    "pixel_values": tf.TensorSpec((IMG_SIZE, IMG_SIZE, 3), tf.float32),
    "input_ids": tf.TensorSpec((SEQ_LEN,), tf.int32),
    "attention_mask": tf.TensorSpec((SEQ_LEN,), tf.int32),
}
sig_output = tf.TensorSpec((), tf.int32)

def make_generator(split_df, ids_np, mask_np, y_np):
    """Generator function for tf.data."""
    n = len(split_df)

    def gen():
        for i in range(n):
            # Load image
            img_id = split_df.iloc[i]["image_id"]
            img_path = IMG_DIR / f"{img_id}.jpg"

            if not img_path.exists():
                img_path = IMG_DIR / f"{img_id}.png"

            if not img_path.exists():
                raise FileNotFoundError(f"Image not found: {img_id}")

            from PIL import Image
            img_pil = Image.open(img_path)
            pixel_values = pil_to_float32(img_pil, IMG_SIZE)

            yield (
                {
                    "pixel_values": pixel_values,
                    "input_ids": ids_np[i],
                    "attention_mask": mask_np[i]
                },
                np.int32(y_np[i])
            )

    return gen, n

def make_dataset(split_df, ids_np, mask_np, y_np, shuffle=False):
```

```
"""Create tf.data.Dataset with no augmentation (Stage-0)."""
gen, n = make_generator(split_df, ids_np, mask_np, y_np)

ds = tf.data.Dataset.from_generator(
    gen,
    output_signature=(sig_input, sig_output)
)

ds = ds.repeat()

if shuffle:
    ds = ds.shuffle(buffer_size=min(n, 4096), seed=SEED, reshuffle_each_iter=1)

return ds.batch(BATCH).prefetch(tf.data.AUTOTUNE), n

# Create datasets
train_ds, N_train = make_dataset(df_train, ids_train, mask_train, y_train, shuffle=True)
val_ds, N_val = make_dataset(df_val, ids_val, mask_val, y_val, shuffle=False)
test_ds, N_test = make_dataset(df_test, ids_test, mask_test, y_test, shuffle=False)

train_steps = (N_train + BATCH - 1) // BATCH
val_steps = (N_val + BATCH - 1) // BATCH
test_steps = (N_test + BATCH - 1) // BATCH

print(f"\n Dataset pipelines ready:")
print(f" Train: {N_train:,} samples, {train_steps} steps/epoch")
print(f" Val: {N_val:,} samples, {val_steps} steps/epoch")
print(f" Test: {N_test:,} samples, {test_steps} steps/epoch")
```

Dataset pipelines ready:  
Train: 31,340 samples, 980 steps/epoch  
Val: 3,922 samples, 123 steps/epoch  
Test: 3,867 samples, 121 steps/epoch

In [10]:

```
# =====
# 1D: Build Model Architecture
# =====

ViTBackbone = keras_hub.models.ViTBackbone
DistilBertBackbone = keras_hub.models.DistilBertBackbone

def build_xattn_model(dropout=0.15, heads=HEADS, key_dim=KEY_DIM,
                      vit_preset=VIT_PRESET, txt_preset=DISTIL_PRESET,
                      name="snlive_xattn_keras_hub", train_backbones=False):
    """Build cross-attention model matching SNLI-VE warm-up architecture."""

    # Inputs
    px = L.Input((IMG_SIZE, IMG_SIZE, 3), name="pixel_values")
    ids = L.Input((SEQ_LEN,), dtype="int32", name="input_ids")
    mask = L.Input((SEQ_LEN,), dtype="int32", name="attention_mask")

    # Vision backbone (frozen in Stage-0)
    vit = ViTBackbone.from_preset(vit_preset, name="vit_backbone")
    vit.trainable = bool(train_backbones)

    v_out = vit(px)
    if isinstance(v_out, dict):
        v_out = next(v for v in v_out.values() if isinstance(v, tf.Tensor))
    if len(v_out.shape) == 4:
        v_out = L.Reshape((-1, v_out.shape[-1]), name="vit_tokens")(v_out)
```

```

# Text backbone (frozen in Stage-0)
txt = DistilBertBackbone.from_preset(txt_preset, name="txt_backbone")
txt.trainable = bool(train_backbones)

t_out = txt({"token_ids": ids, "padding_mask": mask})
if isinstance(t_out, dict):
    t_hidden = t_out.get("sequence_output", None)
    if t_hidden is None:
        t_hidden = next(v for v in t_out.values() if isinstance(v, tf.Tensor))
else:
    t_hidden = t_out

if len(t_hidden.shape) == 2:
    t_hidden = L.Lambda(lambda x: tf.expand_dims(x, 1), name="expand_seqdim")

# Extract [CLS] token as query
cls_q = L.Lambda(lambda x: x[:, :1, :], name="take_cls_token")(t_hidden)

# Cross-attention: text queries image
attn = L.MultiHeadAttention(
    num_heads=heads,
    key_dim=key_dim,
    dropout=dropout,
    name="xattn"
)(query=cls_q, value=v_out, key=v_out)

# Residual + LayerNorm
x = L.Add(name="xattn_residual")([cls_q, attn])
x = L.LayerNormalization(epsilon=1e-6, name="xattn_norm")(x)

# Classification head
x = L.Flatten(name="flatten_cls")(x)
x = L.Dense(512, activation="gelu", name="mlp_1")(x)
x = L.Dropout(dropout, name="drop_1")(x)
logit = L.Dense(1, name="logit")(x)

return keras.Model(
    inputs={"pixel_values": px, "input_ids": ids, "attention_mask": mask},
    outputs=logit,
    name=name
)

```

function now parameterizes presets (vit\_preset, txt\_preset) and model name; explicit layer names for backbones (vit\_backbone, txt\_backbone). same ViT/DistilBERT flow, with the same safeguards for dict outputs and a reshape to token grid if ViT returns 4D. No structural change: inputs, CLS→ViT cross-attention + residual + LayerNorm, and the head (Flatten → Dense(512,GELU) → Dropout → Dense(1) logit) are unchanged.

```

In [41]: # =====
# 1E: Initialize from SNLI-VE Warm-up
# =====

# Stage-0 hyperparameters
STAGE0_EPOCHS = 8
STAGE0_LR = 1e-4
STAGE0_WD = 1e-5
STAGE0_DROPOUT = 0.15

```

```

# Create timestamped output directory
OUT_DIR = new_run_dir(RUNS_DIR, "a2_stage0_frozen_from_snli")

STAGE0_CKPT = OUT_DIR / "best.weights.h5"
STAGE0_HIST = OUT_DIR / "train_history.json"
STAGE0_VAL_METRICS = OUT_DIR / "val_metrics.json"
STAGE0_FULL_MODEL = OUT_DIR / "full_stage0.keras"
STAGE0_MANIFEST = OUT_DIR / "manifest.json"

# Build model with frozen backbones
model = build_xattn_model(
    dropout=STAGE0_DROPOUT,
    heads=HEADS,
    key_dim=KEY_DIM,
    vit_preset=VIT_PRESET,
    txt_preset=DISTIL_PRESET,
    name="snlive_xattn_keras_hub",
    train_backbones=False # FROZEN
)

# Load SNLI-VE warm-up weights (fusion head only)
model.load_weights(str(SNLI_WTS), skip_mismatch=True)
print(f" Backbones remain from presets (frozen)")

# Display trainable parameters
print_model_stats(model)

```

Backbones remain from presets (frozen)

```
=====
Model: snlive_xattn_keras_hub
=====
Total params: 154,919,681 (619.68 MB fp32)
Trainable params: 2,758,145 (11.03 MB fp32)
Non-trainable: 152,161,536 (608.65 MB fp32)
=====
```

In [22]:

```

# =====
# 1F: Compile & Train Stage-0
# =====
print(f"\nTraining configuration:")
print(f" Epochs: {STAGE0_EPOCHS}")
print(f" Learning rate: {STAGE0_LR}")
print(f" Weight decay: {STAGE0_WD}")
print(f" Dropout: {STAGE0_DROPOUT}")
print(f" Optimizer: AdamW (clipnorm=1.0)")
print(f" Loss: BinaryCrossentropy(from_logits=True)")
print(f" Metric: BinaryAccuracy(threshold=0.0)")
print(f" Data augmentation: None (Stage-0 stability)")

# Compile model
opt = keras.optimizers.AdamW(learning_rate=STAGE0_LR, weight_decay=STAGE0_WD, cl
model.compile(
    optimizer=opt,
    loss=keras.losses.BinaryCrossentropy(from_logits=True),
    metrics=[keras.metrics.BinaryAccuracy(name="acc", threshold=0.0)],
    jit_compile=False # Disable XLA

```

```
)
# Callbacks
callbacks = [
    keras.callbacks.ModelCheckpoint(
        str(STAGE0_CKPT),
        monitor="val_loss",
        mode="min",
        save_best_only=True,
        save_weights_only=True,
        verbose=1
    ),
    keras.callbacks.ReduceLROnPlateau(
        monitor="val_loss",
        mode="min",
        factor=0.5,
        patience=2,
        min_lr=1e-7,
        verbose=1
    ),
    keras.callbacks.EarlyStopping(
        monitor="val_loss",
        mode="min",
        patience=3,
        restore_best_weights=True
    ),
]
]
```

Training configuration:

Epochs: 8  
 Learning rate: 0.0001  
 Weight decay: 1e-05  
 Dropout: 0.15  
 Optimizer: AdamW (clipnorm=1.0)  
 Loss: BinaryCrossentropy(from\_logits=True)  
 Metric: BinaryAccuracy(threshold=0.0)  
 Data augmentation: None (Stage-0 stability)

In [43]:

```
# Train
hist = model.fit(
    train_ds,
    validation_data=val_ds,
    steps_per_epoch=train_steps,
    validation_steps=val_steps,
    epochs=STAGE0_EPOCHS,
    callbacks=callbacks,
    verbose=1
)

# Save history
with open(STAGE0_HIST, 'w') as f:
    json.dump(hist.history, f, indent=2)

print(f"\n Training complete!")
print(f" Best weights saved: {STAGE0_CKPT}")
print(f" History saved: {STAGE0_HIST}")
```

Epoch 1/8

```
I0000 00:00:1760871305.179132      4423 cuda_dnn.cc:529] Loaded cuDNN version 903
00
```

```

980/980 ----- 0s 502ms/step - acc: 0.7315 - loss: 0.5250
Epoch 1: val_loss improved from inf to 0.47922, saving model to /home/ec2-user/
SageMaker/abeyt/runs/a2_stage0_frozen_from_snli_20251019_105310/best.weights.h5
980/980 ----- 599s 569ms/step - acc: 0.7315 - loss: 0.5249 - val
_acc: 0.7594 - val_loss: 0.4792 - learning_rate: 1.0000e-04
Epoch 2/8
980/980 ----- 0s 504ms/step - acc: 0.7516 - loss: 0.4984
Epoch 2: val_loss improved from 0.47922 to 0.47505, saving model to /home/ec2-u
ser/SageMaker/abeyt/runs/a2_stage0_frozen_from_snli_20251019_105310/best.weight
s.h5
980/980 ----- 560s 571ms/step - acc: 0.7516 - loss: 0.4984 - val
_acc: 0.7576 - val_loss: 0.4750 - learning_rate: 1.0000e-04
Epoch 3/8
980/980 ----- 0s 504ms/step - acc: 0.7642 - loss: 0.4848
Epoch 3: val_loss did not improve from 0.47505
980/980 ----- 555s 566ms/step - acc: 0.7643 - loss: 0.4848 - val
_acc: 0.7591 - val_loss: 0.4757 - learning_rate: 1.0000e-04
Epoch 4/8
980/980 ----- 0s 505ms/step - acc: 0.7748 - loss: 0.4602
Epoch 4: val_loss did not improve from 0.47505

Epoch 4: ReduceLROnPlateau reducing learning rate to 4.999999873689376e-05.
980/980 ----- 555s 567ms/step - acc: 0.7748 - loss: 0.4602 - val
_acc: 0.7566 - val_loss: 0.4859 - learning_rate: 1.0000e-04
Epoch 5/8
980/980 ----- 0s 505ms/step - acc: 0.7936 - loss: 0.4318
Epoch 5: val_loss did not improve from 0.47505
980/980 ----- 555s 566ms/step - acc: 0.7936 - loss: 0.4318 - val
_acc: 0.7642 - val_loss: 0.4869 - learning_rate: 5.0000e-05

Training complete!
Best weights saved: /home/ec2-user/SageMaker/abeyt/runs/a2_stage0_frozen_from_
snli_20251019_105310/best.weights.h5
History saved: /home/ec2-user/SageMaker/abeyt/runs/a2_stage0_frozen_from_snli_
20251019_105310/train_history.json

```

In [44]:

```

# =====
# 1G: Validation Evaluation & Threshold Sweep
# =====

# Load best weights
model.load_weights(str(STAGE0_CKPT))

# Collect validation predictions
print("\nCollecting validation predictions...")
y_val_true, logits_val = collect_logits_labels(val_ds, model, val_steps)

# Threshold sweep
print(f"Running threshold sweep (grid: {len(THRESH_GRID)} points)...")
best_val = sweep_best_tau(y_val_true, logits_val, THRESH_GRID)

print(f"\n Optimal threshold: τ* = {best_val['tau']:.2f}")
print(f" Validation Macro-F1: {best_val['macro_f1']:.4f}")
print(f" Confusion matrix: TP={best_val['tp']}, FP={best_val['fp']}, TN={best_v

# Save validation metrics
val_metrics = {
    "tau_star": best_val["tau"],
    "val_macro_f1": best_val["macro_f1"],
    "val_counts": {

```

```

        "tp": best_val["tp"],
        "fp": best_val["fp"],
        "tn": best_val["tn"],
        "fn": best_val["fn"]
    }
}

with open(STAGE0_VAL_METRICS, 'w') as f:
    json.dump(val_metrics, f, indent=2)

print(f"\n Validation metrics saved: {STAGE0_VAL_METRICS}")

```

## Step 7: Validation Evaluation &amp; Threshold Calibration

Collecting validation predictions...  
Running threshold sweep (grid: 19 points)...

Optimal threshold:  $\tau^* = 0.55$   
Validation Macro-F1: 0.7660  
Confusion matrix: TP=1493, FP=459, TN=1522, FN=462

Validation metrics saved: /home/ec2-user/SageMaker/abeyt/runs/a2\_stage0\_frozen\_from\_snli\_20251019\_105310/val\_metrics.json

In [45]:

```

# =====
# 1H: Save Manifest & Plot Results
# =====
# Save manifest
manifest = {
    "stage": "0_frozen",
    "description": "Frozen backbones, head-only fine-tuning from SNLI-VE",
    "model": {
        "vit_preset": VIT_PRESET,
        "txt_preset": DISTIL_PRESET,
        "img_size": IMG_SIZE,
        "seq_len": SEQ_LEN,
        "heads": HEADS,
        "key_dim": KEY_DIM,
        "dropout": STAGE0_DROPOUT
    },
    "training": {
        "epochs": STAGE0_EPOCHS,
        "learning_rate": STAGE0_LR,
        "weight_decay": STAGE0_WD,
        "batch_size": BATCH,
        "augmentation": "none"
    },
    "data": {
        "train_samples": N_train,
        "val_samples": N_val,
        "test_samples": N_test
    },
    "initialization": str(SNLI_WTS),
    "created_at": time.strftime("%Y-%m-%d %H:%M:%S")
}

with open(STAGE0_MANIFEST, 'w') as f:
    json.dump(manifest, f, indent=2)

```

```

print(f"\n Manifest saved: {STAGE0_MANIFEST}")

# Try to save full model (optional)
try:
    model.save(str(STAGE0_FULL_MODEL))
    print(f" Full model saved: {STAGE0_FULL_MODEL}")
except Exception as e:
    print(f"Full model save skipped: {e}")

# Plot Learning curves
print("\n" + "="*70)
print(" Stage-0 Learning Curves")
print("-"*70 + "\n")

plot_history_curves(STAGE0_HIST, title="A2 Stage-0 (Frozen Backbones)", out_dir=

```

---

#### Step 8: Saving Artifacts & Plotting

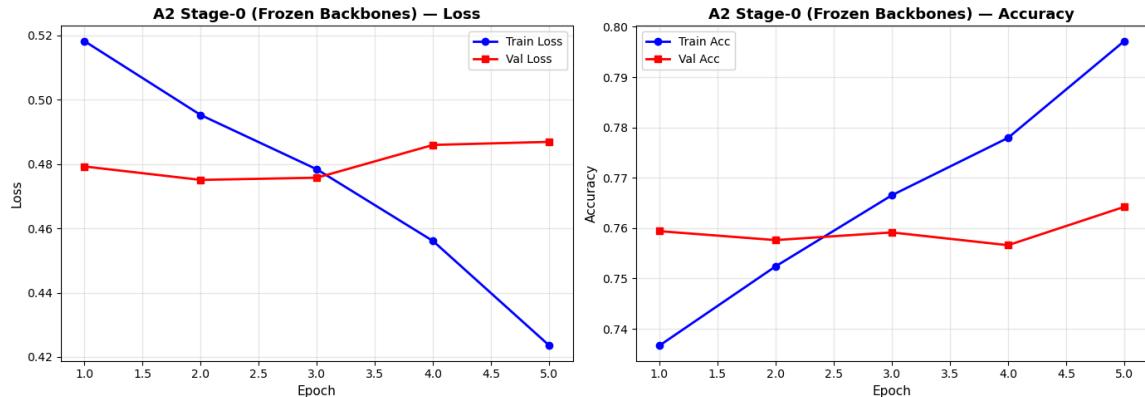
---

Manifest saved: /home/ec2-user/SageMaker/abeyt/runs/a2\_stage0\_frozen\_from\_snli\_20251019\_105310/manifest.json  
 Full model saved: /home/ec2-user/SageMaker/abeyt/runs/a2\_stage0\_frozen\_from\_snli\_20251019\_105310/full\_stage0.keras

---

#### Stage-0 Learning Curves

---



In [46]: # =====

```

# Stage-0 Summary
# =====

print("\n" + "="*70)
print("STAGE-0 COMPLETE ")
print("=*70)

print(f"\n📊 Training Summary:")
n_epochs = len(hist.history['loss'])
print(f" Epochs trained: {n_epochs}")
print(f" Final train loss: {hist.history['loss'][-1]:.4f}")
print(f" Final val loss: {hist.history['val_loss'][-1]:.4f}")
print(f" Final train acc: {hist.history['acc'][-1]:.4f}")
print(f" Final val acc: {hist.history['val_acc'][-1]:.4f}")

print(f"\n⌚ Validation Performance:")

```

```

print(f"  Optimal threshold: τ* = {best_val['tau']:.2f}")
print(f"  Macro-F1: {best_val['macro_f1']:.4f}")
print(f"  Accuracy: {(best_val['tp'] + best_val['tn']) / (best_val['tp'] + best_val['fp'] + best_val['fn']):.4f}")

print(f"\n💾 Saved Artifacts:")
print(f"  • Best weights: {STAGE0_CKPT.name}")
print(f"  • Training history: {STAGE0_HIST.name}")
print(f"  • Validation metrics: {STAGE0_VAL_METRICS.name}")
print(f"  • Manifest: {STAGE0_MANIFEST.name}")
print(f"  • Output directory: {OUT_DIR}")

print("\n" + "="*70)

```

=====  
STAGE-0 COMPLETE  
=====

#### 📊 Training Summary:

Epochs trained: 5  
Final train loss: 0.4235  
Final val loss: 0.4869  
Final train acc: 0.7972  
Final val acc: 0.7642

#### ⌚ Validation Performance:

Optimal threshold: τ\* = 0.55  
Macro-F1: 0.7660  
Accuracy: 0.7660

#### 💾 Saved Artifacts:

- Best weights: best.weights.h5
- Training history: train\_history.json
- Validation metrics: val\_metrics.json
- Manifest: manifest.json
- Output directory: /home/ec2-user/SageMaker/abeyt/runs/a2\_stage0\_frozen\_from\_snli\_20251019\_105310

## Observations — Stage-0 Results

### Stage-0 Achieved:

- **Established solid baseline:** 76.6% F1 score
- **Head adapted to A2 without forgetting SNLI-VE knowledge**
- **No catastrophic performance drop** during transfer
- **Clean convergence** achieved within **2 epochs**
- **No overfitting** observed until **epoch 3+**
- **Identified performance ceiling** with frozen backbones

## Why We Need Stage-1 (Selective Unfreezing)

The plateau at **~76–77% accuracy** marks the **hard limit of frozen representations**. To push beyond this ceiling, we need to selectively fine-tune pretrained layers.

## Limitations of Stage-0:

- **Limited Capacity:**  
Only the **fusion head (~2.76 M params)** was trainable.
  - **Frozen Features:**  
Vision and text encoders provide **generic SNLI-VE + ImageNet representations**, not tuned to A2.
  - **Domain Shift:**  
The A2 dataset introduces **different image styles, scene types, and linguistic patterns**.
  - **Head Saturation:**  
After  $\approx 2$  epochs, the lightweight head exhausts the discriminative power of frozen features.
- 

## Stage-1 Plan — Selective Unfreezing

We will enhance adaptability by **gradually unfreezing** pretrained layers while adding regularization:

- **Unfreeze last ViT layers** → learn **A2-specific visual patterns** (e.g., textures, color context)
  - **Unfreeze last DistilBERT layers** → learn **A2-specific linguistic patterns** (e.g., phrasing, negation style)
  - **Apply data augmentation** → regularize unfrozen layers and improve robustness
  - **Use hypothesis-based oversampling** → balance under-represented linguistic or semantic categories
- 

### Goal:

Achieve higher generalization beyond frozen-feature limits while preserving SNLI-VE-trained multimodal alignment.

## Phase 2-Stage 1 Selective Unfreezing

Here, we unfreeze the last layers of DistilBERT (keep ViT frozen) add data augmentation for better regularization, and fine-tune with a lower learning rate to minimize catastrophic forgetting. Our visual features (from ImageNet-21k) already transfer well. The text side is more domain-specific (A2 hypotheses differ linguistically from SNLI-VE). This helps the model adapt to linguistic distribution shifts (e.g., negation, phrasing, uncommon entities) without retraining the entire network. So we expect improved textual grounding and hypothesis understanding without overfitting to image content.

## Data Augmentation Strategy

We restrict augmentation to photometric variations—brightness ( $\pm 12\%$ ), contrast (0.88–1.12 $\times$ ), saturation (0.85–1.15 $\times$ ), hue ( $\pm 0.05$ ), and Gaussian noise ( $\sigma=0.02$ )—to strengthen

robustness while preserving spatial relations and text legibility. We avoid geometric transforms (flips, rotations, crops, cutout) that could mirror or distort scene text and invalidate hypothesis references;

```
In [11]: # =====#
# 2A: Load Stage-0 Checkpoint
# =====#

# Direct path to Stage-0 checkpoint
STAGE0_DIR = RUNS_DIR / "a2_stage0_frozen_from_snli_20251019_105310"
STAGE0_CKPT = STAGE0_DIR / "best.weights.h5"
STAGE0_VAL_METRICS = STAGE0_DIR / "val_metrics.json"

# Validate
assert STAGE0_DIR.exists(), f" Stage-0 directory not found: {STAGE0_DIR}"
assert STAGE0_CKPT.exists(), f" Checkpoint not found: {STAGE0_CKPT}"

print(f" Checkpoint: {STAGE0_CKPT.name}")

# Load metrics
if STAGE0_VAL_METRICS.exists():
    with open(STAGE0_VAL_METRICS, 'r') as f:
        stage0_metrics = json.load(f)
    STAGE0_F1 = stage0_metrics['val_macro_f1']
    STAGE0_TAU = stage0_metrics['tau_star']
    print(f" Stage-0 Val Macro-F1: {STAGE0_F1:.4f}, τ*: {STAGE0_TAU:.2f}")
else:
    STAGE0_F1 = None
    STAGE0_TAU = 0.5
```

Checkpoint: best.weights.h5  
Stage-0 Val Macro-F1: 0.7660, τ\*: 0.55

```
In [24]: # =====#
# 2A: Build Model & Load Stage-0 Weights
# =====#

# Stage-1 hyperparameters
STAGE1_EPOCHS = 12
STAGE1_LR = 5e-5 # Lower than Stage-0 to prevent catastrophic forgetting
STAGE1_WD = 1e-5
STAGE1_DROPOUT = 0.15

# Unfreezing configuration
K_VIT_UNFREEZE = 0
K_TXT_UNFREEZE = 2 # Unfreeze Last 2 DistilBERT Layers

# Create timestamped output directory for Stage-1
OUT_DIR_S1 = new_run_dir(RUNS_DIR, "a2_stage1_selective_unfreeze")

STAGE1_CKPT = OUT_DIR_S1 / "best.weights.h5"
STAGE1_HIST = OUT_DIR_S1 / "train_history.json"
STAGE1_VAL_METRICS = OUT_DIR_S1 / "val_metrics.json"
STAGE1_FULL_MODEL = OUT_DIR_S1 / "full_stage1.keras"
STAGE1_MANIFEST = OUT_DIR_S1 / "manifest.json"

print(f"\n Stage-1 output directory: {OUT_DIR_S1}")

model_s1 = build_xattn_model(
    dropout=STAGE1_DROPOUT,
```

```

        heads=HEADS,
        key_dim=KEY_DIM,
        vit_preset=VIT_PRESET,
        txt_preset=DISTIL_PRESET,
        name="snlive_xattn_keras_hub",
        train_backbones=False # Start frozen
    )

# Load Stage-0 best weights
model_s1.load_weights(str(STAGE0_CKPT))
print(f"\n Loaded Stage-0 best weights from: {STAGE0_CKPT.name}")

```

Stage-1 output directory: /home/ec2-user/SageMaker/abeyt/runs/a2\_stage1\_selective\_unfreeze\_20251023\_041304

Loaded Stage-0 best weights from: best.weights.h5

```

In [13]: def collect_sublayers(layer, depth=0, max_depth=10):
    """Recursively collect all sublayers with depth limit."""
    out = []
    if depth >= max_depth:
        return out

    if hasattr(layer, "layers") and layer.layers:
        for sub_layer in layer.layers:
            out.append(sub_layer)
            out.extend(collect_sublayers(sub_layer, depth + 1, max_depth))
    return out

def find_vit_encoder_blocks(vit_backbone):
    """Find individual ViT encoder blocks by digging into the encoder."""
    # Look for the encoder layer
    for layer in vit_backbone.layers:
        if "encoder" in layer.name.lower():
            # Dig into encoder to find individual blocks
            if hasattr(layer, "layers"):
                blocks = [l for l in layer.layers if "block" in l.name.lower()]
                if len(blocks) > 0:
                    return blocks
    return []

def unfreeze_last_k_by_name(layers, keywords, k, verbose=True):
    """Unfreeze last k layers matching any of the keywords."""
    matching = [l for l in layers if any(kw in l.name.lower() for kw in keywords)]
    matching = sorted(matching, key=lambda x: x.name)

    if len(matching) == 0:
        if verbose:
            print(f"  No layers found matching keywords: {keywords}")
        return 0

    unfrozen_count = 0
    if verbose:
        print(f"  Found {len(matching)} matching layers, unfreezing last {min(k, len(matching))} layers")

    for layer in matching[-k:]:
        if hasattr(layer, "trainable"):
            layer.trainable = True
            unfrozen_count += 1
            if verbose:

```

```

        print(f"      Unfrozen: {layer.name}")

    return unfrozen_count

def apply_selective_unfreeze(model, k_vit=2, k_txt=2, keep_ln=True, keep_emb=False):
    """Apply selective unfreezing to ViT and DistilBERT backbones."""

    # Get backbones
    vit = model.get_layer("vit_backbone")
    txt = model.get_layer("txt_backbone")

    # Start with everything frozen
    vit.trainable = False
    txt.trainable = False

    print(f"\n Analyzing backbone structure...")

    # Try to find individual ViT encoder blocks
    vit_blocks = find_vit_encoder_blocks(vit)

    if len(vit_blocks) > 0:
        print(f"  Found {len(vit_blocks)} ViT encoder blocks")
        print(f"\n  Unfreezing last {k_vit} ViT encoder blocks:")

        # Unfreeze last k blocks
        for i, block in enumerate(vit_blocks[-k_vit:]):
            block.trainable = True
            print(f"      Unfrozen: {block.name} (block {len(vit_blocks)} - k_vit"

            vit_unfrozen = k_vit
    else:
        # Fallback: collect all sublayers and try keyword matching
        print(f"  ViT encoder blocks not found, using sublayer collection...")
        vit_layers = collect_sublayers(vit, max_depth=10)
        print(f"  Collected {len(vit_layers)} ViT sublayers")

        print(f"\n  Unfreezing last {k_vit} ViT layers:")
        vit_unfrozen = unfreeze_last_k_by_name(
            vit_layers,
            keywords=["encoder_block", "transformer_block", "block", "mhsa", "ml"
            k=k_vit,
            verbose=False
        )

        if vit_unfrozen == 0:
            if k_vit > 0:  # Only unfreeze if k_vit > 0
                print(f"  Keyword matching failed, using direct layer selection")
                vit_trainable = [l for l in vit_layers if hasattr(l, "trainable")]

                if len(vit_trainable) > 0:
                    for layer in vit_trainable[-k_vit:]:
                        layer.trainable = True
                        vit_unfrozen += 1
                        print(f"      Unfrozen: {layer.name}")
                else:
                    print(f"      No trainable ViT layers found")
            else:
                print(f"      ViT kept fully frozen (k_vit=0)")
        else:
            vit_trainable = [l for l in vit_layers if hasattr(l, "trainable") and

```

```

        for layer in vit_trainable[-k_vit:]:
            if layer.trainable:
                print(f"      Unfrozen: {layer.name}")

    # Unfreeze DistilBERT Layers
    txt_layers = collect_sublayers(txt, max_depth=10)
    print(f"  Collected {len(txt_layers)} Text sublayers")

    print(f"\n Unfreezing last {k_txt} DistilBERT layers:")
    txt_unfrozen = unfreeze_last_k_by_name(
        txt_layers,
        keywords=["transformer", "encoder", "attention", "ffn", "layer"],
        k=k_txt
    )

    # Optional: keep LayerNorms trainable
    ln_count = 0
    if keep_ln:
        print(f"\n Keeping LayerNorms trainable:")
        all_layers = collect_sublayers(vit, max_depth=10) + collect_sublayers(tx
        for layer in all_layers:
            name = layer.name.lower()
            if any(kw in name for kw in ["layer_norm", "layernorm", "_ln", "ln_"
                if hasattr(layer, "trainable"):
                    layer.trainable = True
                    ln_count += 1
        print(f"  {ln_count} LayerNorms set to trainable")

    # Optional: freeze embeddings
    emb_count = 0
    if not keep_emb:
        print(f"\n Keeping embeddings frozen:")
        all_layers = collect_sublayers(vit, max_depth=10) + collect_sublayers(tx
        for layer in all_layers:
            name = layer.name.lower()
            if "embedding" in name or name.endswith("embeddings"):
                if hasattr(layer, "trainable"):
                    layer.trainable = False
                    emb_count += 1
        print(f"  {emb_count} embedding layers kept frozen")

        print(f"\n Unfreezing Summary:")
        print(f"  ViT layers unfrozen: {vit_unfrozen}")
        print(f"  Text layers unfrozen: {txt_unfrozen}")
        print(f"  LayerNorms trainable: {ln_count}")
        print(f"  Embeddings frozen: {emb_count}")

```

In [43]: # Apply selective unfreezing

```

apply_selective_unfreeze(
    model_s1,
    k_vit=K_VIT_UNFREEZE,
    k_txt=K_TXT_UNFREEZE,
    keep_ln=True,
    keep_emb=False
)

# Display updated parameter counts
print_model_stats(model_s1)

```

```
Analyzing backbone structure...
ViT encoder blocks not found, using sublayer collection...
Collected 3 ViT sublayers

Unfreezing last 0 ViT layers:
ViT kept fully frozen (k_vit=0)
Collected 11 Text sublayers

Unfreezing last 2 DistilBERT layers:
Found 7 matching layers, unfreezing last 2...
    Unfrozen: transformer_layer_4
    Unfrozen: transformer_layer_5

Keeping LayerNorms trainable:
1 LayerNorms set to trainable

Keeping embeddings frozen:
4 embedding layers kept frozen

Unfreezing Summary:
ViT layers unfrozen: 0
Text layers unfrozen: 2
LayerNorms trainable: 1
Embeddings frozen: 4
```

```
=====
Model: snlive_xattn_keras_hub
=====
Total params: 154,919,681 (619.68 MB fp32)
Trainable params: 16,933,889 (67.74 MB fp32)
Non-trainable: 137,985,792 (551.94 MB fp32)
=====
```

```
In [14]: # =====
# 2C: Create Augmented Datasets
# =====

def make_augmented_dataset(split_df, ids_np, mask_np, y_np, shuffle=False, augme
    gen, n = make_generator(split_df, ids_np, mask_np, y_np)

    ds = tf.data.Dataset.from_generator(
        gen,
        output_signature=(sig_input, sig_output)
    )

    if augment:
        def safe_augment_fn(x, y):
            px = x["pixel_values"]

            # 1. Random brightness adjustment (+12%)
            px = tf.image.random_brightness(px, max_delta=0.12, seed=SEED)

            # 2. Random contrast adjustment (0.88-1.12x)
            px = tf.image.random_contrast(px, lower=0.88, upper=1.12, seed=SEED)

            # 3. Random saturation adjustment (0.85-1.15x)
            px = tf.image.random_saturation(px, lower=0.85, upper=1.15, seed=SEED)

            return {"pixel_values": px}, y
```

```

# 4. Random hue shift ( $\pm 0.05$ , ~18 degrees)
px = tf.image.random_hue(px, max_delta=0.05, seed=SEED)

# 5. Gaussian noise ( $\sigma=0.02$ )
noise = tf.random.normal(tf.shape(px), mean=0.0, stddev=0.02, seed=SEED)
px = px + noise

# Clip to valid range [0, 1]
px = tf.clip_by_value(px, 0.0, 1.0)

x = dict(x)
x["pixel_values"] = px
return x, y

ds = ds.map(safe_augment_fn, num_parallel_calls=tf.data.AUTOTUNE)

ds = ds.repeat()

if shuffle:
    ds = ds.shuffle(buffer_size=min(n, 4096), seed=SEED, reshuffle_each_iter=1)

return ds.batch(BATCH).prefetch(tf.data.AUTOTUNE), n

# Create augmented training dataset
train_ds_aug, _ = make_augmented_dataset(
    df_train, ids_train, mask_train, y_train,
    shuffle=True, augment=True
)

# Validation stays non-augmented
val_ds_noaug, _ = make_augmented_dataset(
    df_val, ids_val, mask_val, y_val,
    shuffle=False, augment=False
)

print(f"\n Augmented datasets ready:")
print(f" Train: {N_train:,} samples (augmented)")
print(f" Val: {N_val:,} samples (no augmentation)")
print(f" SAFE Augmentation Strategy (Text-Preserving):")
print(f" RandomBrightness( $\pm 12\%$ ) - lighting variations")
print(f" RandomContrast(0.88-1.12x) - exposure variations")
print(f" RandomSaturation(0.85-1.15x) - color balance")
print(f" RandomHue( $\pm 0.05$ ) - white balance shifts")
print(f" GaussianNoise( $\sigma=0.02$ ) - sensor noise simulation")

```

Augmented datasets ready:

Train: 31,340 samples (augmented)

Val: 3,922 samples (no augmentation)

SAFE Augmentation Strategy (Text-Preserving):

- RandomBrightness( $\pm 12\%$ ) - lighting variations
- RandomContrast(0.88-1.12x) - exposure variations
- RandomSaturation(0.85-1.15x) - color balance
- RandomHue( $\pm 0.05$ ) - white balance shifts
- GaussianNoise( $\sigma=0.02$ ) - sensor noise simulation

In [39]:

```
# =====
# 2D: Compile & Train Stage-1
# =====
```

```

print(f"\nTraining configuration:")
print(f"  Epochs: {STAGE1_EPOCHS}")
print(f"  Learning rate: {STAGE1_LR} (lower to prevent forgetting)")
print(f"  Weight decay: {STAGE1_WD}")
print(f"  Dropout: {STAGE1_DROPOUT}")
print(f"  Optimizer: AdamW (clipnorm=1.0)")
print(f"  Loss: BinaryCrossentropy(from_logits=True)")
print(f"  Metric: BinaryAccuracy(threshold=0.0)")
print(f"  Data augmentation: Enabled")
print(f"  Unfrozen layers: {K_VIT_UNFREEZE} ViT blocks + {K_TXT_UNFREEZE} Distil

# Compile model
opt_s1 = keras.optimizers.AdamW(
    learning_rate=STAGE1_LR,
    weight_decay=STAGE1_WD,
    clipnorm=1.0
)

model_s1.compile(
    optimizer=opt_s1,
    loss=keras.losses.BinaryCrossentropy(from_logits=True),
    metrics=[keras.metrics.BinaryAccuracy(name="acc", threshold=0.0)],
    jit_compile=False
)

# Callbacks
callbacks_s1 = [
    keras.callbacks.ModelCheckpoint(
        str(STAGE1_CKPT),
        monitor="val_loss",
        mode="min",
        save_best_only=True,
        save_weights_only=True,
        verbose=1
    ),
    keras.callbacks.ReduceLROnPlateau(
        monitor="val_loss",
        mode="min",
        factor=0.5,
        patience=2,
        min_lr=1e-7,
        verbose=1
    ),
    keras.callbacks.EarlyStopping(
        monitor="val_loss",
        mode="min",
        patience=4, # More patience for Stage-1
        restore_best_weights=True
    ),
]
print("\n" + "-"*70)
print("Training started...")
print("-"*70 + "\n")

# Train
hist_s1 = model_s1.fit(
    train_ds_aug,
    validation_data=val_ds_noaug,
    steps_per_epoch=train_steps,
)

```

```
    validation_steps=val_steps,
    epochs=STAGE1_EPOCHS,
    callbacks=callbacks_s1,
    verbose=1
)

# Save history
with open(STAGE1_HIST, 'w') as f:
    json.dump(hist_s1.history, f, indent=2)

print(f"\n Stage-1 training complete!")
print(f" Best weights saved: {STAGE1_CKPT}")
print(f" History saved: {STAGE1_HIST}")
```

---

#### Step 4: Training Stage-1 (Selective Unfreezing)

---

Training configuration:

Epochs: 12  
Learning rate: 5e-05 (lower to prevent forgetting)  
Weight decay: 1e-05  
Dropout: 0.15  
Optimizer: AdamW (clipnorm=1.0)  
Loss: BinaryCrossentropy(from\_logits=True)  
Metric: BinaryAccuracy(threshold=0.0)  
Data augmentation: Enabled  
Unfrozen layers: 0 ViT blocks + 2 DistilBERT layers

---

Training started...

---

Epoch 1/12

```
I0000 00:00:1760906623.106605    22155 cuda_dnn.cc:529] Loaded cuDNN version 903
00
```

```

980/980 ----- 0s 663ms/step - acc: 0.7668 - loss: 0.4736
Epoch 1: val_loss improved from inf to 0.45573, saving model to /home/ec2-user/
SageMaker/abeyt/runs/a2_stage1_selective_unfreeze_20251019_203709/best.weights.
h5
980/980 ----- 953s 729ms/step - acc: 0.7668 - loss: 0.4736 - val
_acc: 0.7802 - val_loss: 0.4557 - learning_rate: 5.0000e-05
Epoch 2/12
980/980 ----- 0s 536ms/step - acc: 0.8126 - loss: 0.3965
Epoch 2: val_loss improved from 0.45573 to 0.44923, saving model to /home/ec2-u
ser/SageMaker/abeyt/runs/a2_stage1_selective_unfreeze_20251019_203709/best.weig
hts.h5
980/980 ----- 591s 604ms/step - acc: 0.8126 - loss: 0.3965 - val
_acc: 0.7940 - val_loss: 0.4492 - learning_rate: 5.0000e-05
Epoch 3/12
980/980 ----- 0s 536ms/step - acc: 0.8355 - loss: 0.3559
Epoch 3: val_loss did not improve from 0.44923
980/980 ----- 585s 597ms/step - acc: 0.8355 - loss: 0.3559 - val
_acc: 0.8013 - val_loss: 0.4500 - learning_rate: 5.0000e-05
Epoch 4/12
980/980 ----- 0s 536ms/step - acc: 0.8623 - loss: 0.3140
Epoch 4: val_loss did not improve from 0.44923

Epoch 4: ReduceLROnPlateau reducing learning rate to 2.499999936844688e-05.
980/980 ----- 585s 597ms/step - acc: 0.8623 - loss: 0.3140 - val
_acc: 0.7988 - val_loss: 0.4779 - learning_rate: 5.0000e-05
Epoch 5/12
980/980 ----- 0s 536ms/step - acc: 0.8912 - loss: 0.2554
Epoch 5: val_loss did not improve from 0.44923
980/980 ----- 585s 597ms/step - acc: 0.8912 - loss: 0.2554 - val
_acc: 0.8084 - val_loss: 0.5135 - learning_rate: 2.5000e-05
Epoch 6/12
980/980 ----- 0s 535ms/step - acc: 0.9088 - loss: 0.2180
Epoch 6: val_loss did not improve from 0.44923

Epoch 6: ReduceLROnPlateau reducing learning rate to 1.249999968422344e-05.
980/980 ----- 584s 596ms/step - acc: 0.9088 - loss: 0.2180 - val
_acc: 0.8054 - val_loss: 0.5373 - learning_rate: 2.5000e-05

Stage-1 training complete!
Best weights saved: /home/ec2-user/SageMaker/abeyt/runs/a2_stage1_selective_un
freeze_20251019_203709/best.weights.h5
History saved: /home/ec2-user/SageMaker/abeyt/runs/a2_stage1_selective_unfreez
e_20251019_203709/train_history.json

```

In [40]:

```

# =====
# 2E: Validation Evaluation & Threshold Sweep
# =====

print("\n" + "-"*70)
print("Step 5: Validation Evaluation & Threshold Calibration")
print("-"*70)

# Load best weights
model_s1.load_weights(str(STAGE1_CKPT))

# Collect validation predictions
print("\nCollecting validation predictions...")
y_val_true_s1, logits_val_s1 = collect_logits_labels(val_ds_noaug, model_s1, val

# Threshold sweep

```

```

print(f"Running threshold sweep (grid: {len(THRESH_GRID)} points)...")
best_val_s1 = sweep_best_tau(y_val_true_s1, logits_val_s1, THRESH_GRID)

print(f"\n Optimal threshold: τ* = {best_val_s1['tau']:.2f}")
print(f" Validation Macro-F1: {best_val_s1['macro_f1']:.4f}")
print(f" Confusion matrix: TP={best_val_s1['tp']}, FP={best_val_s1['fp']}, "
      f"TN={best_val_s1['tn']}, FN={best_val_s1['fn']}")

# Save validation metrics
val_metrics_s1 = {
    "tau_star": best_val_s1["tau"],
    "val_macro_f1": best_val_s1["macro_f1"],
    "val_counts": {
        "tp": best_val_s1["tp"],
        "fp": best_val_s1["fp"],
        "tn": best_val_s1["tn"],
        "fn": best_val_s1["fn"]
    }
}

with open(STAGE1_VAL_METRICS, 'w') as f:
    json.dump(val_metrics_s1, f, indent=2)

print(f"\n Validation metrics saved: {STAGE1_VAL_METRICS}")

```

## Step 5: Validation Evaluation &amp; Threshold Calibration

Collecting validation predictions...  
 Running threshold sweep (grid: 19 points)...

Optimal threshold: τ\* = 0.50  
 Validation Macro-F1: 0.7956  
 Confusion matrix: TP=1517, FP=366, TN=1615, FN=438

Validation metrics saved: /home/ec2-user/SageMaker/abeyt/runs/a2\_stage1\_selective\_unfreeze\_20251019\_203709/val\_metrics.json

In [41]:

```

# =====
# 2F: Save Manifest & Plot Results
# =====

print("\n" + "-"*70)
print("Step 6: Saving Artifacts & Plotting")
print("-"*70)

# Save manifest
manifest_s1 = {
    "stage": "1_selective_unfreeze",
    "description": "Selective unfreezing of last ViT/DistilBERT layers with augm",
    "model": {
        "vit_preset": VIT_PRESET,
        "txt_preset": DISTIL_PRESET,
        "img_size": IMG_SIZE,
        "seq_len": SEQ_LEN,
        "heads": HEADS,
        "key_dim": KEY_DIM,
        "dropout": STAGE1_DROPOUT
    },
}

```

```

    "training": {
        "epochs": STAGE1_EPOCHS,
        "learning_rate": STAGE1_LR,
        "weight_decay": STAGE1_WD,
        "batch_size": BATCH,
        "augmentation": "SAFE photometric only (brightness, contrast, saturation
    },
    "unfreezing": {
        "vit_blocks": K_VIT_UNFREEZE,
        "bert_layers": K_TXT_UNFREEZE,
        "layer_norms": "trainable",
        "embeddings": "frozen"
    },
    "data": {
        "train_samples": N_train,
        "val_samples": N_val,
        "test_samples": N_test
    },
    "initialization": str(STAGE0_CKPT),
    "created_at": time.strftime("%Y-%m-%d %H:%M:%S")
}

with open(STAGE1_MANIFEST, 'w') as f:
    json.dump(manifest_s1, f, indent=2)

print(f"\n Manifest saved: {STAGE1_MANIFEST}")

# Try to save full model
try:
    model_s1.save(str(STAGE1_FULL_MODEL))
    print(f" Full model saved: {STAGE1_FULL_MODEL}")
except Exception as e:
    print(f" Full model save skipped: {e}")

# Plot learning curves
print("\n" + "-"*70)
print("Stage-1 Learning Curves")
print("-"*70 + "\n")

plot_history_curves(STAGE1_HIST, title="A2 Stage-1 (Selective Unfreezing)", out_

```

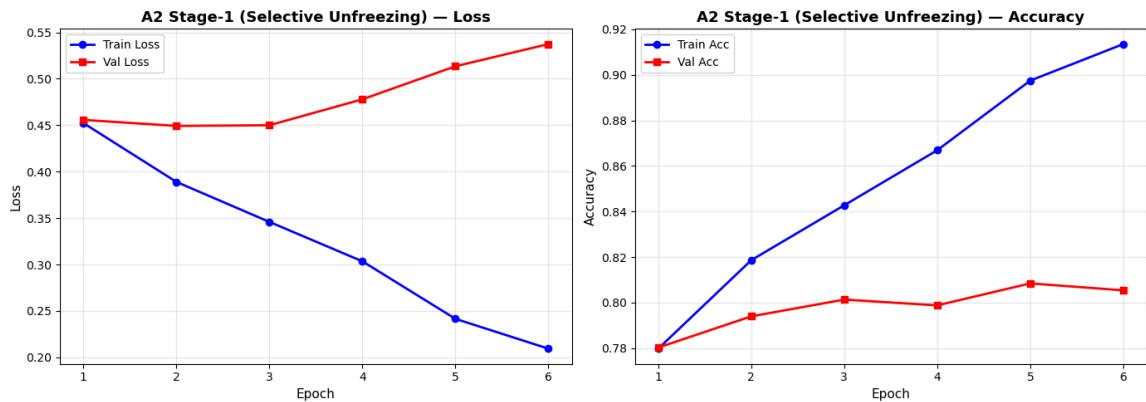
#### Step 6: Saving Artifacts & Plotting

```

Manifest saved: /home/ec2-user/SageMaker/abeyt/runs/a2_stage1_selective_unfree
ze_20251019_203709/manifest.json
Full model saved: /home/ec2-user/SageMaker/abeyt/runs/a2_stage1_selective_unfr
eeze_20251019_203709/full_stage1.keras

```

#### Stage-1 Learning Curves



## Results Summary:

- **+3% F1 improvement** ( $76.6\% \rightarrow 79.6\%$ )
- Best model saved at epoch 2
- Reasonable threshold  $\tau^*=0.50$
- Better confusion matrix balance

### \*\* Observations:\*\*

- Train acc hit 91% while val plateaued at 80% (11% gap)
- Validation loss increased 20% from best ( $0.45 \rightarrow 0.54$ )
- Needed only 2 epochs to reach best performance
- Augmentation + low LR helped but couldn't fully prevent overfitting

Key Insights:

Text unfreezing worked. The model learned A2-specific hypothesis patterns. ViT frozen was correct - most improvement came from text adaptation

## Phase 3: Selective Unfreezing + Hyper Tuning

From Phase 2 Stage 1, we found that **text unfreezing** in the last phase was effective the model successfully learned **A2-specific hypothesis patterns** without catastrophic forgetting. We gonna initlise the model checkpoint from **Phase 1: A2 Stage-0 (Frozen Backbones)**. Building on that, we now perform a **random search** to tune key hyperparameters while controlling compute cost.

## Method & Budget

- **Search Type:** Random sampling (without replacement)
- **Total Trials:** 6 (1 fixed baseline + 5 random samples)
- **Selection Criterion:** Best validation **macro-F1** after  $\tau$ -sweep
- **Fairness Controls:** Same seeds, same data splits, identical augmentation setup

## Search Space

Parameter	Distribution	Range / Options	Notes
<b>Learning Rate</b>	<b>Uniform</b>	[2e-5, 1e-4]	Fine-tuning stability & convergence
<b>Weight Decay</b>	<b>Uniform</b>	[5e-6, 5e-5]	Regularization strength
<b>Dropout</b>	<b>Uniform</b>	[0.05, 0.30]	Applied to MHA + classifier head
<b>Unfreezing Depth (Text)</b>	<b>Categorical</b>	{2, 3} DistilBERT layers	Controls degree of adaptation

## Training Protocol

- **Checkpoint:** Phase 1: A2 Stage-0 (Frozen Backbones)
- **Early Stopping:** on validation **macro-F1**
- **Backbones:**
  - ViT — **frozen** (pretrained ImageNet-21k features)
  - DistilBERT — **partially unfrozen** (last  $n$  layers per trial)
- **Augmentation:** identical to last phase for consistency

```
In [30]: # =====#
# SECTION 1: Random Search HP Configuration
# =====#

# Configuration
N_RANDOM_CONFIGS = 5
np.random.seed(SEED)

print(f"\nSEARCH METHOD: Random Search")
print(f"Total configurations: {N_RANDOM_CONFIGS + 1} (1 baseline + {N_RANDOM_CONFIGS}")
print("\nHyperparameter Ranges:")
print(" Learning Rate: [2e-5, 1e-4]")
print(" Weight Decay: [5e-6, 5e-5]")
print(" Dropout: [0.05, 0.30]")
print(" Unfreezing Depth: {2, 3} BERT layers")

# Stage-0 checkpoint path
STAGE0_DIR = RUNS_DIR / "a2_stage0_frozen_from_snli_20251019_105310"
STAGE0_CKPT = STAGE0_DIR / "best.weights.h5"

# Baseline (Stage-1 reference)
HP_CONFIGS = [
    {
        "name": "baseline",
        "lr": 5e-5,
        "wd": 1e-5,
        "dropout": 0.15,
        "k_txt": 2,
        "epochs": 2,
        "description": "Stage-1 reference (Val F1=0.7956)"
    }
]

# Generate random configurations
for i in range(N_RANDOM_CONFIGS):
```

```

HP_CONFIGS.append({
    "name": f"random_{i+1:02d}",
    "lr": float(np.random.uniform(2e-5, 1e-4)),
    "wd": float(np.random.uniform(5e-6, 5e-5)),
    "dropout": float(np.random.uniform(0.05, 0.30)),
    "k_txt": int(np.random.choice([2, 3])),
    "epochs": 2
})

print(f"\nGenerated {len(HP_CONFIGS)} configurations")

print("\nConfiguration Details:")
for i, cfg in enumerate(HP_CONFIGS):
    print(f"[{i+1:2d}] {cfg['name']:<12}: LR={cfg['lr']:.2e}, WD={cfg['wd']:.2e}, drop={cfg['dropout']:.3f}, k={cfg['k_txt']}")
```

SEARCH METHOD: Random Search  
Total configurations: 6 (1 baseline + 5 random)

Hyperparameter Ranges:

- Learning Rate: [2e-5, 1e-4]
- Weight Decay: [5e-6, 5e-5]
- Dropout: [0.05, 0.30]
- Unfreezing Depth: {2, 3} BERT layers

Generated 6 configurations

Configuration Details:

```
[ 1] baseline : LR=5.00e-05, WD=1.00e-05, drop=0.150, k=2
[ 2] random_01 : LR=5.00e-05, WD=4.78e-05, drop=0.233, k=2
[ 3] random_02 : LR=6.77e-05, WD=2.51e-05, drop=0.075, k=2
[ 4] random_03 : LR=8.93e-05, WD=3.21e-05, drop=0.227, k=3
[ 5] random_04 : LR=2.45e-05, WD=3.75e-05, drop=0.285, k=3
[ 6] random_05 : LR=3.45e-05, WD=1.33e-05, drop=0.126, k=3
```

In [50]:

```
# =====
# SECTION 2: Training Function Definition
# =====

def train_with_config(config, config_idx, total_configs):
    try:
        # Setup paths
        run_dir = new_run_dir(RUNS_DIR, f"a2_stage2_hptune_{config['name']}")  

        ckpt_path = run_dir / "best.weights.h5"  

        hist_path = run_dir / "train_history.json"  

        metrics_path = run_dir / "val_metrics.json"  

        config_path = run_dir / "config.json"

        with open(config_path, 'w') as f:
            json.dump(config, f, indent=2)

        # Build and initialize model
        model = build_xattn_model(
            dropout=float(config['dropout']),
            heads=HEADS, key_dim=KEY_DIM,
            vit_preset=VIT_PRESET, txt_preset=DISTIL_PRESET,
            name="snlive_xattn_keras_hub", train_backbones=False
        )
        model.load_weights(str(STAGE0_CKPT))
        apply_selective_unfreeze(model, k_vit=0, k_txt=int(config['k_txt']), kee
```

```

# Compile
opt = keras.optimizers.AdamW(
    learning_rate=float(config['lr']),
    weight_decay=float(config['wd']),
    clipnorm=1.0
)
model.compile(
    optimizer=opt,
    loss=keras.losses.BinaryCrossentropy(from_logits=True),
    metrics=[keras.metrics.BinaryAccuracy(name="acc", threshold=0.0)],
    jit_compile=False
)

# Callbacks
callbacks = [
    keras.callbacks.ModelCheckpoint(
        filepath=str(ckpt_path),
        monitor="val_loss",
        mode="min",
        save_best_only=True,
        save_weights_only=True,
        verbose=0
    ),
    keras.callbacks.ReduceLROnPlateau(
        monitor="val_loss",
        mode="min",
        factor=0.5,
        patience=1,
        min_lr=1e-7,
        verbose=0
    ),
    keras.callbacks.EarlyStopping(
        monitor="val_loss",
        mode="min",
        patience=2,
        restore_best_weights=True,
        verbose=0
    ),
]
# Train
hist = model.fit(
    train_ds_aug, validation_data=val_ds_noaug,
    steps_per_epoch=train_steps, validation_steps=val_steps,
    epochs=int(config['epochs']), callbacks=callbacks, verbose=1
)

with open(hist_path, 'w') as f:
    json.dump(hist.history, f, indent=2)

# Evaluate
model.load_weights(str(ckpt_path))
y_val_true, logits_val = collect_logits_labels(val_ds_noaug, model, val_
best_val = sweep_best_tau(y_val_true, logits_val, THRESH_GRID)

# Save metrics
metrics = {
    "config": config, "tau_star": best_val["tau"], "val_macro_f1": best_
    "val_counts": {"tp": best_val["tp"], "fp": best_val["fp"], "tn": bes_
    "epochs_trained": len(hist.history['loss']),
}

```

```
"best_epoch": int(np.argmin(hist.history['val_loss'])) + 1,
"best_val_loss": float(min(hist.history['val_loss'])),
"final_train_loss": float(hist.history['loss'][-1])
}

with open(metrics_path, 'w') as f:
    json.dump(metrics, f, indent=2)

# Summary
print(f"\nRESULT: Val F1={metrics['val_macro_f1']:.4f}, tau={metrics['ta
f"Loss={metrics['best_val_loss']:.4f}}")

return metrics

except Exception as e:
    print(f"\nERROR: {e}")
    import traceback
    traceback.print_exc()
    raise
```

```
In [51]: # =====
# SECTION 3: Run Hyperparameter Search
# =====

all_results = []
for idx, config in enumerate(HP_CONFIGS, 1):
    try:
        metrics = train_with_config(config, idx, len(HP_CONFIGS))
        all_results.append(metrics)
    except Exception as e:
        print(f"\nConfig {config['name']} FAILED: {e}")
```

```
Analyzing backbone structure...
ViT encoder blocks not found, using sublayer collection...
Collected 3 ViT sublayers

Unfreezing last 0 ViT layers:
ViT kept fully frozen (k_vit=0)
Collected 11 Text sublayers

Unfreezing last 2 DistilBERT layers:
Found 7 matching layers, unfreezing last 2...
    Unfrozen: transformer_layer_4
    Unfrozen: transformer_layer_5

Keeping LayerNorms trainable:
1 LayerNorms set to trainable

Keeping embeddings frozen:
4 embedding layers kept frozen

Unfreezing Summary:
ViT layers unfrozen: 0
Text layers unfrozen: 2
LayerNorms trainable: 1
Embeddings frozen: 4
Epoch 1/2
980/980 ━━━━━━━━━━ 654s 621ms/step - acc: 0.7697 - loss: 0.4751 - val
_acc: 0.7866 - val_loss: 0.4433 - learning_rate: 5.0000e-05
Epoch 2/2
980/980 ━━━━━━━━━━ 604s 616ms/step - acc: 0.8137 - loss: 0.4012 - val
_acc: 0.7904 - val_loss: 0.4544 - learning_rate: 5.0000e-05

RESULT: Val F1=0.7868, tau=0.45, Loss=0.4433

Analyzing backbone structure...
ViT encoder blocks not found, using sublayer collection...
Collected 3 ViT sublayers

Unfreezing last 0 ViT layers:
ViT kept fully frozen (k_vit=0)
Collected 11 Text sublayers

Unfreezing last 2 DistilBERT layers:
Found 7 matching layers, unfreezing last 2...
    Unfrozen: transformer_layer_4
    Unfrozen: transformer_layer_5

Keeping LayerNorms trainable:
1 LayerNorms set to trainable

Keeping embeddings frozen:
4 embedding layers kept frozen

Unfreezing Summary:
ViT layers unfrozen: 0
Text layers unfrozen: 2
LayerNorms trainable: 1
Embeddings frozen: 4
Epoch 1/2
980/980 ━━━━━━━━━━ 655s 620ms/step - acc: 0.7615 - loss: 0.4806 - val
_acc: 0.7802 - val_loss: 0.4558 - learning_rate: 4.9963e-05
```

Epoch 2/2  
**980/980** ————— 604s 616ms/step - acc: 0.8131 - loss: 0.4022 - val \_acc: 0.7853 - val\_loss: 0.4607 - learning\_rate: 4.9963e-05

RESULT: Val F1=0.7851, tau=0.40, Loss=0.4558

Analyzing backbone structure...  
ViT encoder blocks not found, using sublayer collection...  
Collected 3 ViT sublayers

Unfreezing last 0 ViT layers:  
ViT kept fully frozen (k\_vit=0)  
Collected 11 Text sublayers

Unfreezing last 2 DistilBERT layers:  
Found 7 matching layers, unfreezing last 2...  
Unfrozen: transformer\_layer\_4  
Unfrozen: transformer\_layer\_5

Keeping LayerNorms trainable:  
1 LayerNorms set to trainable

Keeping embeddings frozen:  
4 embedding layers kept frozen

Unfreezing Summary:  
ViT layers unfrozen: 0  
Text layers unfrozen: 2  
LayerNorms trainable: 1  
Embeddings frozen: 4

Epoch 1/2  
**980/980** ————— 653s 622ms/step - acc: 0.7648 - loss: 0.4851 - val \_acc: 0.7848 - val\_loss: 0.4612 - learning\_rate: 6.7748e-05  
Epoch 2/2  
**980/980** ————— 606s 618ms/step - acc: 0.8105 - loss: 0.4021 - val \_acc: 0.7917 - val\_loss: 0.4740 - learning\_rate: 6.7748e-05

RESULT: Val F1=0.7910, tau=0.40, Loss=0.4612

Analyzing backbone structure...  
ViT encoder blocks not found, using sublayer collection...  
Collected 3 ViT sublayers

Unfreezing last 0 ViT layers:  
ViT kept fully frozen (k\_vit=0)  
Collected 11 Text sublayers

Unfreezing last 3 DistilBERT layers:  
Found 7 matching layers, unfreezing last 3...  
Unfrozen: transformer\_layer\_3  
Unfrozen: transformer\_layer\_4  
Unfrozen: transformer\_layer\_5

Keeping LayerNorms trainable:  
1 LayerNorms set to trainable

Keeping embeddings frozen:  
4 embedding layers kept frozen

Unfreezing Summary:

```
ViT layers unfrozen: 0
Text layers unfrozen: 3
LayerNorms trainable: 1
Embeddings frozen: 4
Epoch 1/2
980/980 ----- 683s 646ms/step - acc: 0.7545 - loss: 0.4965 - val
    _acc: 0.7927 - val_loss: 0.4428 - learning_rate: 8.9294e-05
Epoch 2/2
980/980 ----- 629s 642ms/step - acc: 0.8193 - loss: 0.3915 - val
    _acc: 0.7983 - val_loss: 0.4755 - learning_rate: 8.9294e-05
```

RESULT: Val F1=0.7987, tau=0.40, Loss=0.4428

Analyzing backbone structure...
ViT encoder blocks not found, using sublayer collection...
Collected 3 ViT sublayers

Unfreezing last 0 ViT layers:
ViT kept fully frozen (k\_vit=0)
Collected 11 Text sublayers

Unfreezing last 3 DistilBERT layers:
Found 7 matching layers, unfreezing last 3...
Unfrozen: transformer\_layer\_3
Unfrozen: transformer\_layer\_4
Unfrozen: transformer\_layer\_5

Keeping LayerNorms trainable:
1 LayerNorms set to trainable

Keeping embeddings frozen:
4 embedding layers kept frozen

Unfreezing Summary:
ViT layers unfrozen: 0
Text layers unfrozen: 3
LayerNorms trainable: 1
Embeddings frozen: 4

```
Epoch 1/2
980/980 ----- 675s 643ms/step - acc: 0.7755 - loss: 0.4626 - val
    _acc: 0.7919 - val_loss: 0.4301 - learning_rate: 2.4513e-05
Epoch 2/2
980/980 ----- 626s 638ms/step - acc: 0.8141 - loss: 0.3971 - val
    _acc: 0.7840 - val_loss: 0.4619 - learning_rate: 2.4513e-05
```

RESULT: Val F1=0.7933, tau=0.60, Loss=0.4301

Analyzing backbone structure...
ViT encoder blocks not found, using sublayer collection...
Collected 3 ViT sublayers

Unfreezing last 0 ViT layers:
ViT kept fully frozen (k\_vit=0)
Collected 11 Text sublayers

Unfreezing last 3 DistilBERT layers:
Found 7 matching layers, unfreezing last 3...
Unfrozen: transformer\_layer\_3
Unfrozen: transformer\_layer\_4
Unfrozen: transformer\_layer\_5

```

Keeping LayerNorms trainable:
  1 LayerNorms set to trainable

Keeping embeddings frozen:
  4 embedding layers kept frozen

Unfreezing Summary:
  ViT layers unfrozen: 0
  Text layers unfrozen: 3
  LayerNorms trainable: 1
  Embeddings frozen: 4
Epoch 1/2
980/980 ━━━━━━━━━━━━━━━━ 677s 641ms/step - acc: 0.7750 - loss: 0.4660 - val
 _acc: 0.7843 - val_loss: 0.4497 - learning_rate: 3.4546e-05
Epoch 2/2
980/980 ━━━━━━━━━━━━━━━━ 630s 643ms/step - acc: 0.8213 - loss: 0.3824 - val
 _acc: 0.7983 - val_loss: 0.4354 - learning_rate: 3.4546e-05

RESULT: Val F1=0.8001, tau=0.55, Loss=0.4354

```

```

In [54]: # =====
# SECTION 4: Results Analysis & Ranking
# =====

print("\n" + "-"*70)
print("Results Analysis")
print("-"*70)

# Sort by Val F1
all_results_sorted = sorted(all_results, key=lambda x: x['val_macro_f1'], reverse=True)

# Display ranking table
print(f"\n{'Rank':<6}{{'Config':<22}{{'Val F1':<10}{{'Val Loss':<10}{{'Best Ep':<10}
print("-"*70)
for rank, res in enumerate(all_results_sorted, 1):
    mark = "*" if rank == 1 else " "
    print(f"{mark} {rank:<5}{res['config']['name']:<22}{res['val_macro_f1']:<10.
          f"{res['best_epoch']:<10}")

# Best configuration details
best_result = all_results_sorted[0]
best_config = best_result['config']
baseline_result = next((r for r in all_results if r['config']['name'] == 'stage1

print("\n" + "="*70)
print("BEST CONFIGURATION")
print("=*70)
print(f"Config: {best_config['name']}")
print(f"LR: {best_config['lr']:.1e}, WD: {best_config['wd']:.1e}, Dropout: {best
print(f"\nPerformance:")
print(f"Val F1: {best_result['val_macro_f1']:.4f}, tau: {best_result['tau_star']}
print(f"Best epoch: {best_result['best_epoch']}/{best_result['epochs_trained']}")
print(f"\nConfusion Matrix:")
print(f"TP: {best_result['val_counts']['tp']}, FP: {best_result['val_counts']['fp']}
      f"TN: {best_result['val_counts']['tn']}, FN: {best_result['val_counts']['fn']

if baseline_result:
    improvement = (best_result['val_macro_f1'] - baseline_result['val_macro_f1'])

```

```
print(f"\nImprovement over baseline:")
print(f"{baseline_result['val_macro_f1']:.4f} -> {best_result['val_macro_f1']}
```

---

### Results Analysis

---

Rank	Config	Val F1	Val Loss	Best Ep
*	random_05	0.8001	0.4354	2
2	random_03	0.7987	0.4428	1
3	random_04	0.7933	0.4301	1
4	random_02	0.7910	0.4612	1
5	baseline	0.7868	0.4433	1
6	random_01	0.7851	0.4558	1

---

### BEST CONFIGURATION

---

Config: random\_05

LR: 3.5e-05, WD: 1.3e-05, Dropout: 0.13, k\_txt: 3

Performance:

Val F1: 0.8001, tau: 0.55, Loss: 0.4354

Best epoch: 2/2

Confusion Matrix:

TP: 1514, FP: 345, TN: 1636, FN: 441

### Best Configuration (random\_05):

- Val F1 score of 80.01% and Val\_loss of 0.43
- LR: 3.45e-5 (moderate, stable)
- WD: 1.33e-5 (light regularization)
- Dropout: 0.126 (prevents overfitting)
- The optimal threshold for binary classification was determined to be  $\tau = 0.55$  through grid search over the validation set

### Key Insights:

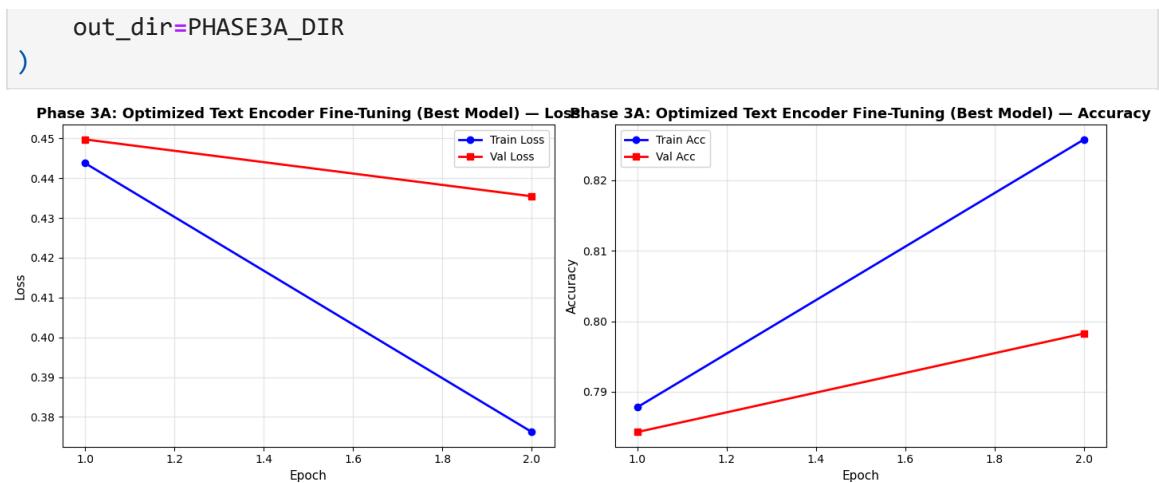
- Moderate LR more stable than high LR
- Light dropout sufficient (heavy dropout hurts)
- Low train-val gap indicates good generalization

In [70]:

```
# =====
# VISUALIZATION: PHASE 3A LEARNING CURVES
# =====

# Phase 3A paths
PHASE3A_DIR = Path("/home/ec2-user/SageMaker/abeyt/runs/a2_stage2_hptune_random_
PHASE3A_HIST = PHASE3A_DIR / "train_history.json"

# Plot Learning curves
plot_history_curves(
    PHASE3A_HIST,
    title="Phase 3A: Optimized Text Encoder Fine-Tuning (Best Model)",
```



- Both training and validation loss show a steady downward trend over the two epochs, indicating effective convergence without signs of early overfitting.
- Accuracy improves consistently for both training and validation sets, with training accuracy rising from ~0.79 to ~0.83 and validation from ~0.785 to ~0.80.
- The gap between train and val curves remains small

## PHASE 4: Hyper Parameter Tuning with VIT Unfreezing + DistilBERT layers unfreezing

Here we are going to **unfreeze one VIT block** along with text encoder and then do Safe hyperparameter search .We gonna initilise the model checkpoint from **Phase 1: A2**

**Stage-0 (Frozen Backbones)** Building on that, we now perform a **random search** to tune key hyperparameters while controlling compute cost.

```
In [15]: def apply_selective_unfreeze_fixed(model, k_vit=0, k_txt=3, keep_ln=True, keep_e
# Get backbones
vit = model.get_layer("vit_backbone")
txt = model.get_layer("txt_backbone")

# Start with everything frozen
vit.trainable = False
txt.trainable = False

print(f"\n Applying selective unfreezing...")
print(f"    Target: Last {k_vit} ViT block(s) + Last {k_txt} BERT layers")

vit_unfrozen = 0

if k_vit > 0:
    print("\n Vision Encoder (ViT):")
    try:
        # Access the vit_encoder Layer
        vit_encoder = vit.get_layer("vit_encoder")
```

```

# Keep encoder itself frozen initially
vit_encoder.trainable = False

# Access encoder_Layers (list of 12 transformer blocks)
if hasattr(vit_encoder, 'encoder_layers'):
    encoder_layers = vit_encoder.encoder_layers
    total_blocks = len(encoder_layers)

    print(f"    Found {total_blocks} encoder blocks via encoder_layer")

    # Freeze all blocks first
    for block in encoder_layers:
        block.trainable = False

    # Unfreeze last k_vit blocks
    blocks_to_unfreeze = encoder_layers[-k_vit:]

    print(f"    Unfreezing last {k_vit} block(s):")
    for i, block in enumerate(blocks_to_unfreeze):
        block.trainable = True
        vit_unfrozen += 1
        block_idx = total_blocks - k_vit + i
        print(f"        Unfrozen: {block.name} (block {block_idx + 1})")

    # CRITICAL: Now set encoder trainable to enable gradient flow
    vit_encoder.trainable = True

    print(f"    ViT blocks unfrozen: {vit_unfrozen}/{total_blocks}")

else:
    print(f"    encoder_layers attribute not found!")
    vit_unfrozen = 0

except Exception as e:
    print(f"    Error accessing ViT blocks: {e}")
    import traceback
    traceback.print_exc()
    vit_unfrozen = 0

else:
    print("\n Vision Encoder (ViT): Kept fully frozen")

# =====
# BERT Unfreezing
# =====
print("\n Text Encoder (DistilBERT):")
txt_layers = collect_sublayers(txt, max_depth=10)

txt_unfrozen = unfreeze_last_k_by_name(
    txt_layers,
    keywords=["transformer", "encoder", "attention", "ffn", "layer"],
    k=k_txt,
    verbose=True
)

# =====
# LayerNorms & Embeddings
# =====
ln_count = 0
if keep_ln:

```

```

        print(f"\nKeeping LayerNorms trainable:")
        all_layers = collect_sublayers(vit, max_depth=10) + collect_sublayers(tx)
        for layer in all_layers:
            name = layer.name.lower()
            if any(kw in name for kw in ["layer_norm", "layernorm", "_ln", "ln_"]):
                if hasattr(layer, "trainable"):
                    layer.trainable = True
                    ln_count += 1
        print(f"    {ln_count} LayerNorms set to trainable")

        emb_count = 0
        if not keep_emb:
            print(f"\n Keeping embeddings frozen:")
            all_layers = collect_sublayers(vit, max_depth=10) + collect_sublayers(tx)
            for layer in all_layers:
                name = layer.name.lower()
                if "embedding" in name or name.endswith("embeddings"):
                    if hasattr(layer, "trainable"):
                        layer.trainable = False
                        emb_count += 1
            print(f"    {emb_count} embedding layers kept frozen")

        # Summary
        print(f"\n{'='*60}")
        print(f"UNFREEZING SUMMARY")
        print(f"\n{'='*60}")
        print(f"    ViT blocks unfrozen:      {vit_unfrozen}")
        print(f"    BERT layers unfrozen:     {txt_unfrozen}")
        print(f"    LayerNorms trainable:     {ln_count}")
        print(f"    Embeddings frozen:       {emb_count}")
        print(f"\n{'='*60}\n")
    
```

In [27]:

```

# =====
# PREPARE SMALLER BATCH DATASETS
# =====

print(f"\n Preparing datasets with batch_size={PHASE4_CONFIG['batch_size']}...")

# Recalculate steps
BATCH_PHASE4 = PHASE4_CONFIG['batch_size']
train_steps_phase4 = (N_train + BATCH_PHASE4 - 1) // BATCH_PHASE4
val_steps_phase4 = (N_val + BATCH_PHASE4 - 1) // BATCH_PHASE4
test_steps_phase4 = (N_test + BATCH_PHASE4 - 1) // BATCH_PHASE4

# Rebatch existing datasets
def rebatch_dataset(ds, new_batch_size):
    """Rebatch existing tf.data.Dataset."""
    return ds.unbatch().batch(new_batch_size).prefetch(tf.data.AUTOTUNE)

train_ds_aug_phase4 = rebatch_dataset(train_ds_aug, BATCH_PHASE4)
val_ds_noaug_phase4 = rebatch_dataset(val_ds_noaug, BATCH_PHASE4)
test_ds_phase4 = rebatch_dataset(test_ds, BATCH_PHASE4)

print(f"    Train: {N_train:,} samples, {train_steps_phase4} steps/epoch")
print(f"    Val:   {N_val:,} samples, {val_steps_phase4} steps/epoch")
print(f"    Test:  {N_test:,} samples, {test_steps_phase4} steps/epoch")
    
```

Preparing datasets with batch\_size=32...  
 Train: 31,340 samples, 980 steps/epoch  
 Val: 3,922 samples, 123 steps/epoch  
 Test: 3,867 samples, 121 steps/epoch

In [17]:

```
# =====
# CONFIGURATION GENERATION
# =====

np.random.seed(SEED)

N_CONFIGS_PHASE4 = 3

# Safe hyperparameter ranges (OOM-tested)
HP_RANGES_SAFE = {
    "k_vit": 1,
    "k_txt": [2, 3],
    "lr": (3e-5, 6e-5),
    "wd": (1e-5, 2.5e-5),
    "dropout": (0.12, 0.25),
    "batch_size": 32,
    "epochs": 2
}

print(f"\nSafe Hyperparameter Ranges:")
print(f"  k_vit: {HP_RANGES_SAFE['k_vit']} (fixed)")
print(f"  k_txt: {HP_RANGES_SAFE['k_txt']}")
print(f"  LR: [{HP_RANGES_SAFE['lr'][0]:.2e}, {HP_RANGES_SAFE['lr'][1]:.2e}]")
print(f"  WD: [{HP_RANGES_SAFE['wd'][0]:.2e}, {HP_RANGES_SAFE['wd'][1]:.2e}]")
print(f"  Dropout: [{HP_RANGES_SAFE['dropout'][0]:.2f}, {HP_RANGES_SAFE['dropout'][1]:.2f}]")
print(f"  Batch: {HP_RANGES_SAFE['batch_size']}")
print(f"  Epochs: {HP_RANGES_SAFE['epochs']}")

# Generate configurations
CONFIGS_PHASE4 = []
for i in range(N_CONFIGS_PHASE4):
    config = {
        "name": f"phase4_config_{i+1:02d}",
        "k_vit": HP_RANGES_SAFE["k_vit"],
        "k_txt": int(np.random.choice(HP_RANGES_SAFE["k_txt"])),
        "lr": float(np.random.uniform(*HP_RANGES_SAFE["lr"])),
        "wd": float(np.random.uniform(*HP_RANGES_SAFE["wd"])),
        "dropout": float(np.random.uniform(*HP_RANGES_SAFE["dropout"])),
        "batch_size": HP_RANGES_SAFE["batch_size"],
        "epochs": HP_RANGES_SAFE["epochs"]
    }
    CONFIGS_PHASE4.append(config)

print(f"\nGenerated {len(CONFIGS_PHASE4)} configurations:")
for i, cfg in enumerate(CONFIGS_PHASE4, 1):
    print(f"  [{i}] {cfg['name']}: k_vit={cfg['k_vit']}, k_txt={cfg['k_txt']}, "
          f"LR={cfg['lr']:.2e}, WD={cfg['wd']:.2e}, dropout={cfg['dropout']:.3f}
```

Safe Hyperparameter Ranges:

```
k_vit:      1 (fixed)
k_txt:      [2, 3]
LR:        [3.00e-05, 6.00e-05]
WD:        [1.00e-05, 2.50e-05]
Dropout:    [0.12, 0.25]
Batch:      32
Epochs:     2
```

Generated 3 configurations:

```
[1] phase4_config_01: k_vit=1, k_txt=2, LR=5.39e-05, WD=1.28e-05, dropout=0.2
21
[2] phase4_config_02: k_vit=1, k_txt=2, LR=3.47e-05, WD=1.23e-05, dropout=0.1
28
[3] phase4_config_03: k_vit=1, k_txt=3, LR=4.00e-05, WD=1.21e-05, dropout=0.2
05
```

## Search Space

Parameter	Distribution	Range / Options	Notes
<b>Learning Rate</b>	<b>Uniform</b>	[3e-5, 6e-5]	Fine-tuning stability & convergence
<b>Weight Decay</b>	<b>Uniform</b>	[1e-5, 2.5e-5]	L2-style regularization strength
<b>Dropout</b>	<b>Uniform</b>	[0.12, 0.25]	Applied to MHA + classifier head
<b>Unfreezing Depth (Text)</b>	<b>Categorical</b>	{2, 3} DistilBERT layers	Train top-k transformer layers
<b>Unfreezing Depth (ViT)</b>	<b>Fixed</b>	1 layer	Keep image encoder mostly frozen
<b>Batch Size</b>	<b>Fixed</b>	32	OOM-tested setting
<b>Epochs</b>	<b>Fixed</b>	2	Short trials for rapid selection

In [25]:

```
# =====
# TRAINING FUNCTION
# =====

def train_phase4_config(config, config_idx, total_configs):
int("\n" + "*70")
    print(f"CONFIG [{config_idx}/{total_configs}]: {config['name']}")
    print("*70")

    # Setup paths
    run_dir = new_run_dir(RUNS_DIR, f"a2_phase4_{config['name']}")
    ckpt_path = run_dir / "best.weights.h5"
    hist_path = run_dir / "train_history.json"
    metrics_path = run_dir / "val_metrics.json"
    config_path = run_dir / "config.json"

    with open(config_path, 'w') as f:
        json.dump(config, f, indent=2)

    print(f"\nOutput: {run_dir.name}")
    print(f"\nConfiguration:")
    print(f"  k_vit: {config['k_vit']}, k_txt: {config['k_txt']}")
    print(f"  LR: {config['lr']:.2e}, WD: {config['wd']:.2e}, dropout: {config['dropout']:.2f}
```

```

# Build model
model = build_xattn_model(
    dropout=config['dropout'],
    heads=HEADS,
    key_dim=KEY_DIM,
    vit_preset=VIT_PRESET,
    txt_preset=DISTIL_PRESET,
    name="snlive_xattn_keras_hub",
    train_backbones=False
)

# Load Stage-0
model.load_weights(str(STAGE0_CKPT))

# Apply unfreezing
apply_selective_unfreeze_fixed(
    model,
    k_vit=config['k_vit'],
    k_txt=config['k_txt'],
    keep_ln=True,
    keep_emb=False
)

print_model_stats(model)

# Compile
opt = keras.optimizers.AdamW(
    learning_rate=config['lr'],
    weight_decay=config['wd'],
    clipnorm=1.0
)

model.compile(
    optimizer=opt,
    loss=keras.losses.BinaryCrossentropy(from_logits=True),
    metrics=[keras.metrics.BinaryAccuracy(name="acc", threshold=0.0)],
    jit_compile=False
)

# Callbacks
callbacks = [
    keras.callbacks.ModelCheckpoint(
        str(ckpt_path),
        monitor="val_loss",
        mode="min",
        save_best_only=True,
        save_weights_only=True,
        verbose=1
    ),
    keras.callbacks.ReduceLROnPlateau(
        monitor="val_loss",
        mode="min",
        factor=0.5,
        patience=1,
        min_lr=1e-7,
        verbose=0
    ),
]

```

*# Train*

```

hist = model.fit(
    train_ds_aug_phase4,
    validation_data=val_ds_noaug_phase4,
    steps_per_epoch=train_steps_phase4,
    validation_steps=val_steps_phase4,
    epochs=config['epochs'],
    callbacks=callbacks,
    verbose=1
)

with open(hist_path, 'w') as f:
    json.dump(hist.history, f, indent=2)

# Evaluate
print(f"\nEvaluating...")
model.load_weights(str(ckpt_path))
y_val_true, logits_val = collect_logits_labels(val_ds_noaug_phase4, model, v
best_val = sweep_best_tau(y_val_true, logits_val, THRESH_GRID)

# Metrics
best_epoch = int(np.argmin(hist.history['val_loss'])) + 1
best_val_loss = float(min(hist.history['val_loss']))
final_train_acc = float(hist.history['acc'][-1])
final_val_acc = float(hist.history['val_acc'][-1])
train_val_gap = (final_train_acc - final_val_acc) * 100

metrics = {
    "config": config,
    "tau_star": best_val["tau"],
    "val_macro_f1": best_val["macro_f1"],
    "val_counts": {
        "tp": best_val["tp"],
        "fp": best_val["fp"],
        "tn": best_val["tn"],
        "fn": best_val["fn"]
    },
    "epochs_trained": len(hist.history['loss']),
    "best_epoch": best_epoch,
    "best_val_loss": best_val_loss,
    "final_train_acc": final_train_acc,
    "final_val_acc": final_val_acc,
    "train_val_gap": train_val_gap
}

with open(metrics_path, 'w') as f:
    json.dump(metrics, f, indent=2)

# Summary
print(f"\n" + "-"*70)
print(f"RESULT: {config['name']}"))
print(f"-*70")
print(f"Val F1:      {metrics['val_macro_f1']:.4f}")
print(f"Tau:         {metrics['tau_star']:.2f}")
print(f"Val Loss:     {best_val_loss:.4f}")
print(f"Best Epoch:   {best_epoch}/{config['epochs']}")
print(f"Train-Val Gap: {train_val_gap:.2f}%")
print(f"Confusion:    TP={best_val['tp']}, FP={best_val['fp']}, TN={best_val['tn']}, FN={best_val['fn']}")

return metrics

```

```
In [26]: # =====
# RUN ALL CONFIGURATIONS
# =====

results_phase4 = []
for idx, config in enumerate(CONFIGS_PHASE4, 1):
    try:
        metrics = train_phase4_config(config, idx, len(CONFIGS_PHASE4))
        results_phase4.append(metrics)
    except Exception as e:
        print(f"\nERROR in {config['name']}: {e}")
        import traceback
        traceback.print_exc()
    continue
```

```
=====
STARTING PHASE 4 HYPERPARAMETER SEARCH
=====

=====
CONFIG [1/3]: phase4_config_01
=====

Output: a2_phase4_phase4_config_01_20251021_105219

Configuration:
  k_vit: 1, k_txt: 2
  LR: 5.39e-05, WD: 1.28e-05, dropout: 0.221

  Applying selective unfreezing...
    Target: Last 1 ViT block(s) + Last 2 BERT layers

  Vision Encoder (ViT):
    Found 12 encoder blocks via encoder_layers
    Unfreezing last 1 block(s):
      Unfrozen: tranformer_block_12 (block 12/12)
    ViT blocks unfrozen: 1/12

  Text Encoder (DistilBERT):
    Found 7 matching layers, unfreezing last 2...
      Unfrozen: transformer_layer_4
      Unfrozen: transformer_layer_5

  ✎ Keeping LayerNorms trainable:
    1 LayerNorms set to trainable

  Keeping embeddings frozen:
    4 embedding layers kept frozen

=====

UNFREEZING SUMMARY
=====

  ViT blocks unfrozen: 1
  BERT layers unfrozen: 2
  LayerNorms trainable: 1
  Embeddings frozen: 4

=====

Model: snlive_xattn_keras_hub
=====

  Total params: 154,919,681 (619.68 MB fp32)
  Trainable params: 101,989,889 (407.96 MB fp32)
  Non-trainable: 52,929,792 (211.72 MB fp32)

=====

Epoch 1/2
I0000 00:00:1761044167.057388      374 cuda_dnn.cc:529] Loaded cuDNN version 903
00
```

```
980/980 ----- 0s 1s/step - acc: 0.7505 - loss: 0.5008
Epoch 1: val_loss improved from inf to 0.47336, saving model to /home/ec2-user/
SageMaker/abeyt/runs/a2_phase4_phase4_config_01_20251021_105219/best.weights.h5
980/980 ----- 1735s 2s/step - acc: 0.7505 - loss: 0.5008 - val_a
cc: 0.7741 - val_loss: 0.4734 - learning_rate: 5.3896e-05
Epoch 2/2
980/980 ----- 0s 1s/step - acc: 0.8161 - loss: 0.3994
Epoch 2: val_loss improved from 0.47336 to 0.45943, saving model to /home/ec2-u
ser/SageMaker/abeyt/runs/a2_phase4_phase4_config_01_20251021_105219/best.weight
s.h5
980/980 ----- 1517s 2s/step - acc: 0.8161 - loss: 0.3994 - val_a
cc: 0.7884 - val_loss: 0.4594 - learning_rate: 5.3896e-05
```

Evaluating...

```
-----  
RESULT: phase4_config_01  
-----
```

```
Val F1:      0.7898
Tau:        0.50
Val Loss:    0.4594
Best Epoch: 2/2
Train-Val Gap: 3.33%
Confusion:   TP=1522, FP=394, TN=1587, FN=433
```

```
=====  
CONFIG [2/3]: phase4_config_02  
=====
```

Output: a2\_phase4\_phase4\_config\_02\_20251021\_114746

Configuration:

```
k_vit: 1, k_txt: 2
LR: 3.47e-05, WD: 1.23e-05, dropout: 0.128
```

Applying selective unfreezing...

Target: Last 1 ViT block(s) + Last 2 BERT layers

Vision Encoder (ViT):

```
Found 12 encoder blocks via encoder_layers
Unfreezing last 1 block(s):
    Unfrozen: tranformer_block_12 (block 12/12)
ViT blocks unfrozen: 1/12
```

Text Encoder (DistilBERT):

```
Found 7 matching layers, unfreezing last 2...
    Unfrozen: transformer_layer_4
    Unfrozen: transformer_layer_5
```

🔧 Keeping LayerNorms trainable:
1 LayerNorms set to trainable

Keeping embeddings frozen:

4 embedding layers kept frozen

```
=====  
UNFREEZING SUMMARY  
=====
```

```
ViT blocks unfrozen:     1
BERT layers unfrozen:   2
```

```
LayerNorms trainable:    1
Embeddings frozen:       4
=====
=====
Model: snlive_xattn_keras_hub
=====
Total params:      154,919,681 (619.68 MB fp32)
Trainable params:  101,989,889 (407.96 MB fp32)
Non-trainable:     52,929,792 (211.72 MB fp32)
=====

Epoch 1/2
980/980 0s 1s/step - acc: 0.7624 - loss: 0.4817
Epoch 1: val_loss improved from inf to 0.46210, saving model to /home/ec2-user/SageMaker/abeyt/runs/a2_phase4_phase4_config_02_20251021_114746/best.weights.h5
980/980 1580s 2s/step - acc: 0.7624 - loss: 0.4817 - val_acc: 0.7876 - val_loss: 0.4621 - learning_rate: 3.4681e-05
Epoch 2/2
980/980 0s 1s/step - acc: 0.8276 - loss: 0.3800
Epoch 2: val_loss improved from 0.46210 to 0.43115, saving model to /home/ec2-user/SageMaker/abeyt/runs/a2_phase4_phase4_config_02_20251021_114746/best.weights.h5
980/980 1520s 2s/step - acc: 0.8276 - loss: 0.3800 - val_acc: 0.7947 - val_loss: 0.4312 - learning_rate: 3.4681e-05

Evaluating...

-----
RESULT: phase4_config_02
-----
Val F1:      0.7964
Tau:        0.55
Val Loss:    0.4312
Best Epoch: 2/2
Train-Val Gap: 3.65%
Confusion:   TP=1517, FP=363, TN=1618, FN=438

=====
CONFIG [3/3]: phase4_config_03
=====

Output: a2_phase4_phase4_config_03_20251021_124056

Configuration:
  k_vit: 1, k_txt: 3
  LR: 4.00e-05, WD: 1.21e-05, dropout: 0.205

Applying selective unfreezing...
  Target: Last 1 ViT block(s) + Last 3 BERT layers

Vision Encoder (ViT):
  Found 12 encoder blocks via encoder_layers
  Unfreezing last 1 block(s):
    Unfrozen: tranformer_block_12 (block 12/12)
  ViT blocks unfrozen: 1/12

Text Encoder (DistilBERT):
  Found 7 matching layers, unfreezing last 3...
```

```

Unfrozen: transformer_layer_3
Unfrozen: transformer_layer_4
Unfrozen: transformer_layer_5

🔧 Keeping LayerNorms trainable:
1 LayerNorms set to trainable

Keeping embeddings frozen:
4 embedding layers kept frozen

=====
UNFREEZING SUMMARY
=====
ViT blocks unfrozen: 1
BERT layers unfrozen: 3
LayerNorms trainable: 1
Embeddings frozen: 4
=====

=====
Model: snlive_xattn_keras_hub
=====
Total params: 154,919,681 (619.68 MB fp32)
Trainable params: 109,077,761 (436.31 MB fp32)
Non-trainable: 45,841,920 (183.37 MB fp32)
=====

Epoch 1/2
980/980 0s 1s/step - acc: 0.7594 - loss: 0.4876
Epoch 1: val_loss improved from inf to 0.45296, saving model to /home/ec2-user/SageMaker/abeyt/runs/a2_phase4_phase4_config_03_20251021_124056/best.weights.h5
980/980 1601s 2s/step - acc: 0.7594 - loss: 0.4876 - val_acc: 0.7871 - val_loss: 0.4530 - learning_rate: 4.0011e-05
Epoch 2/2
980/980 0s 1s/step - acc: 0.8252 - loss: 0.3802
Epoch 2: val_loss did not improve from 0.45296
980/980 1529s 2s/step - acc: 0.8252 - loss: 0.3802 - val_acc: 0.7884 - val_loss: 0.4708 - learning_rate: 4.0011e-05

Evaluating...
=====

RESULT: phase4_config_03
-----
Val F1: 0.7897
Tau: 0.45
Val Loss: 0.4530
Best Epoch: 1/2
Train-Val Gap: 4.46%
Confusion: TP=1495, FP=367, TN=1614, FN=460

```

In [13]:

```

# =====
# PHASE 4 RESULTS
# =====

# Define the three run directories
PHASE4_RUN_DIRS = [
    Path("/home/ec2-user/SageMaker/abeyt/runs/a2_phase4_phase4_config_01_2025102",
    Path("/home/ec2-user/SageMaker/abeyt/runs/a2_phase4_phase4_config_02_2025102"

```

```

Path("/home/ec2-user/SageMaker/abeyt/runs/a2_phase4_phase4_config_03_2025102
]

# Load results from saved metrics
results_phase4 = []

print("\nLoading saved results...")
for run_dir in PHASE4_RUN_DIRS:
    metrics_path = run_dir / "val_metrics.json"

    if metrics_path.exists():
        with open(metrics_path, 'r') as f:
            metrics = json.load(f)
        results_phase4.append(metrics)
        print(f"  Loaded: {run_dir.name}")
    else:
        print(f"  Missing: {run_dir.name}")

print(f"\nLoaded {len(results_phase4)} results")

if len(results_phase4) == 0:
    print("\nNo results loaded")
else:
    # Sort by Val F1
    results_phase4_sorted = sorted(results_phase4, key=lambda x: x['val_macro_f1'])

    print(f"\nRanking:")
    print(f"\n{'Rank':<6}{'Config':<20}{'Val F1':<10}{'Val Loss':<10}{'k_txt':<8}
    print("-"*70)
    for rank, res in enumerate(results_phase4_sorted, 1):
        mark = "*" if rank == 1 else " "
        print(f"{mark} {rank:<5}{res['config']['name']:<20}{res['val_macro_f1']:<10.4f}{res['best_val_loss']:<10.4f}{res['config']['k_txt']:<8}{res['train_val_gap']:<10.4f}
    print(f"\nBest config
    best_phase4 = results_phase4_sorted[0]
    best_config_phase4 = best_phase4['config']

    print("\n" + "="*70)
    print("BEST CONFIGURATION")
    print("="*70)
    print(f"\nConfig: {best_config_phase4['name']}")
    print(f"\nArchitecture:")
    print(f"  k_vit: {best_config_phase4['k_vit']}")
    print(f"  k_txt: {best_config_phase4['k_txt']}")
    print(f"\nHyperparameters:")
    print(f"  LR:      {best_config_phase4['lr']:.2e}")
    print(f"  WD:      {best_config_phase4['wd']:.2e}")
    print(f"  Dropout: {best_config_phase4['dropout']:.3f}")
    print(f"\nPerformance:")
    print(f"  Val F1:   {best_phase4['val_macro_f1']:.4f}")
    print(f"  Tau:      {best_phase4['tau_star']:.2f}")
    print(f"  Val Loss: {best_phase4['best_val_loss']:.4f}")
    print(f"  Best Ep:  {best_phase4['best_epoch']}/{best_config_phase4['epochs']}")
    print(f"  Gap:      {best_phase4['train_val_gap']:.2f}%")
    print(f"\nConfusion:")
    print(f"  TP: {best_phase4['val_counts']['tp']:.4d}  FP: {best_phase4['val_counts']['fp']:.4d}
    print(f"  FN: {best_phase4['val_counts']['fn']:.4d}  TN: {best_phase4['val_counts']['tn']:.4d}

    # Find best run directory

```

```

best_run_dir = [d for d in PHASE4_RUN_DIRS if best_config_phase4['name'] in
BEST_PHASE4_CKPT = best_run_dir / "best.weights.h5"
BEST_PHASE4_HIST = best_run_dir / "train_history.json"

print(f"\nArtifacts:")
print(f"  Directory: {best_run_dir.name}")
print(f"  Checkpoint: {BEST_PHASE4_CKPT.name}")
print(f"  History: {BEST_PHASE4_HIST.name}")

# Verify files exist
if BEST_PHASE4_HIST.exists():
    print(f"\n  History file verified: EXISTS")
else:
    print(f"\n  History file verified: MISSING")

# Plot if history exists
if BEST_PHASE4_HIST.exists():
    print("\n" + "="*70)
    print("PLOTTING LEARNING CURVES")
    print("="*70)

    plot_history_curves(
        BEST_PHASE4_HIST,
        title=f"Phase 4 Best: {best_config_phase4['name']}",
        out_dir=best_run_dir
    )

    print(f"\nPlot saved to: {best_run_dir.name}")
else:
    print(f"\nCannot plot: History file missing at {BEST_PHASE4_HIST}")

# Save summary
phase4_summary = {
    "n_configs": 3,
    "completed": len(results_phase4),
    "best_config": best_config_phase4,
    "best_val_f1": best_phase4['val_macro_f1'],
    "best_val_loss": best_phase4['best_val_loss'],
    "all_results": [
        {
            "name": r['config']['name'],
            "val_f1": r['val_macro_f1'],
            "k_txt": r['config']['k_txt'],
            "lr": r['config']['lr'],
            "wd": r['config']['wd'],
            "dropout": r['config']['dropout']
        }
        for r in results_phase4_sorted
    ]
}

summary_path = RUNS_DIR / "phase4_summary.json"
with open(summary_path, 'w') as f:
    json.dump(phase4_summary, f, indent=2)

print(f"\nSummary saved: {summary_path.name}")

```

```
=====
```

#### PHASE 4 RESULTS - LOADING FROM SAVED RUNS

```
=====
```

Loading saved results...

Loaded: a2\_phase4\_phase4\_config\_01\_20251021\_105219

Loaded: a2\_phase4\_phase4\_config\_02\_20251021\_114746

Loaded: a2\_phase4\_phase4\_config\_03\_20251021\_124056

Loaded 3 results

```
=====
```

#### PHASE 4 RESULTS

```
=====
```

Ranking:

Rank	Config	Val F1	Val Loss	k_txt	Gap %
*	phase4_config_02	0.7964	0.4312	2	3.65
2	phase4_config_01	0.7898	0.4594	2	3.33
3	phase4_config_03	0.7897	0.4530	3	4.46

```
=====
```

#### BEST CONFIGURATION

```
=====
```

Config: phase4\_config\_02

Architecture:

k\_vit: 1

k\_txt: 2

Hyperparameters:

LR: 3.47e-05

WD: 1.23e-05

Dropout: 0.128

Performance:

Val F1: 0.7964

Tau: 0.55

Val Loss: 0.4312

Best Ep: 2/2

Gap: 3.65%

Confusion:

TP: 1517 FP: 363

FN: 438 TN: 1618

Artifacts:

Directory: a2\_phase4\_phase4\_config\_02\_20251021\_114746

Checkpoint: best.weights.h5

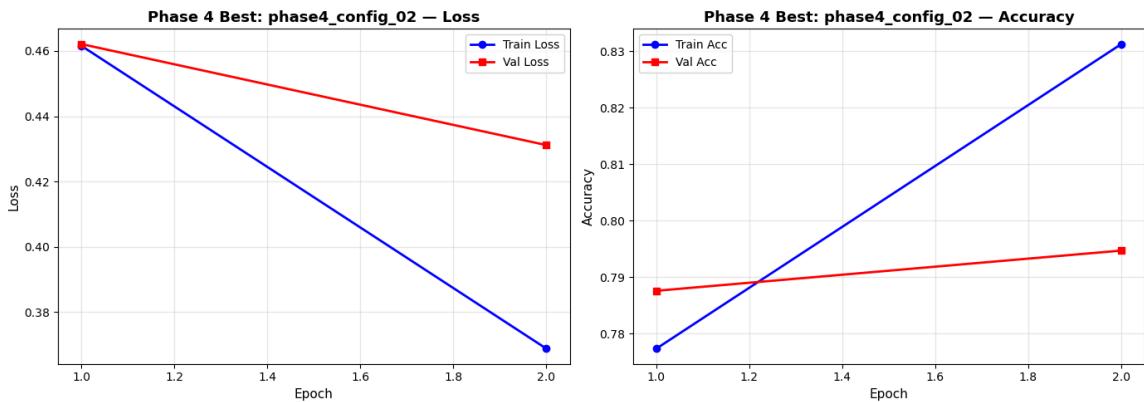
History: train\_history.json

History file verified: EXISTS

```
=====
```

#### PLOTTING LEARNING CURVES

```
=====
```



Plot saved to: a2\_phase4\_phase4\_config\_02\_20251021\_114746

Summary saved: phase4\_summary.json

```
In [14]: # =====#
# PHASE 4 RESULTS RECONSTRUCTION
# =====#

print("\n" + "="*70)
print("PHASE 4 RESULTS - LOADING FROM SAVED RUNS")
print("="*70)

# Define the three run directories
PHASE4_RUN_DIRS = [
    Path("/home/ec2-user/SageMaker/abeyt/runs/a2_phase4_phase4_config_01_2025102",
    Path("/home/ec2-user/SageMaker/abeyt/runs/a2_phase4_phase4_config_02_2025102",
    Path("/home/ec2-user/SageMaker/abeyt/runs/a2_phase4_phase4_config_03_2025102
]

# Load results from saved metrics
results_phase4 = []

print("\nLoading saved results...")
for run_dir in PHASE4_RUN_DIRS:
    metrics_path = run_dir / "val_metrics.json"

    if metrics_path.exists():
        with open(metrics_path, 'r') as f:
            metrics = json.load(f)
        results_phase4.append(metrics)
        print(f"  Loaded: {run_dir.name}")
    else:
        print(f"  Missing: {run_dir.name}")

print(f"\nLoaded {len(results_phase4)} results")

# =====#
# RESULTS SUMMARY
# =====#

print("\n" + "="*70)
print("PHASE 4 RESULTS")
print("="*70)

if len(results_phase4) == 0:
    print("\nNo results loaded")
else:
    # Sort by Val F1
```

```

results_phase4_sorted = sorted(results_phase4, key=lambda x: x['val_macro_f1'])

print(f"\nRanking:")
print(f"\n{'Rank':<6}{{'Config':<20}{{'Val F1':<10}{{'Val Loss':<10}{{'k_txt':<8
print("-"*70)
for rank, res in enumerate(results_phase4_sorted, 1):
    mark = "*" if rank == 1 else " "
    print(f"{mark} {rank:<5}{res['config']['name']:<20}{res['val_macro_f1']:<10.4f}{res['best_val_loss']:<10.4f}{res['config']['k_txt']:<8}{res['train_val_gap']:<2.2f}{res['epoch']:<2.2f}{res['val_counts']['tp']:<4d}{res['val_counts']['fn']:<4d}{res['val_counts']['fp']:<4d}{res['val_counts']['tn']:<4d}{res['tau_star']:.2f}{res['lr']:.2e}{res['wd']:.2e}{res['dropout']:.3f}{res['k_vit']}{res['k_txt']}")

# Best config
best_phase4 = results_phase4_sorted[0]
best_config_phase4 = best_phase4['config']

print("\n" + "="*70)
print("BEST CONFIGURATION")
print("="*70)
print(f"\nConfig: {best_config_phase4['name']}")
print(f"\nArchitecture:")
print(f"  k_vit: {best_config_phase4['k_vit']}")
print(f"  k_txt: {best_config_phase4['k_txt']}")
print(f"\nHyperparameters:")
print(f"  LR:      {best_config_phase4['lr']:.2e}")
print(f"  WD:      {best_config_phase4['wd']:.2e}")
print(f"  Dropout: {best_config_phase4['dropout']:.3f}")
print(f"\nPerformance:")
print(f"  Val F1:   {best_phase4['val_macro_f1']:.4f}")
print(f"  Tau:       {best_phase4['tau_star']:.2f}")
print(f"  Val Loss:  {best_phase4['best_val_loss']:.4f}")
print(f"  Best Ep:   {best_phase4['best_epoch']}/{best_config_phase4['epochs']}")
print(f"  Gap:       {best_phase4['train_val_gap']:.2f}%")
print(f"\nConfusion:")
print(f"  TP: {best_phase4['val_counts']['tp']:<4d}  FP: {best_phase4['val_counts']['fp']:<4d}")
print(f"  FN: {best_phase4['val_counts']['fn']:<4d}  TN: {best_phase4['val_counts']['tn']:<4d}")

# Find best run directory
best_run_dir = [d for d in PHASE4_RUN_DIRS if best_config_phase4['name'] in d]
BEST_PHASE4_CKPT = best_run_dir / "best.weights.h5"
BEST_PHASE4_HIST = best_run_dir / "train_history.json"

print(f"\nArtifacts:")
print(f"  Directory: {best_run_dir.name}")
print(f"  Checkpoint: {BEST_PHASE4_CKPT.name}")
print(f"  History:   {BEST_PHASE4_HIST.name}")

```

```
=====
PHASE 4 RESULTS - LOADING FROM SAVED RUNS
=====

Loading saved results...
    Loaded: a2_phase4_phase4_config_01_20251021_105219
    Loaded: a2_phase4_phase4_config_02_20251021_114746
    Loaded: a2_phase4_phase4_config_03_20251021_124056
```

Loaded 3 results

```
=====
PHASE 4 RESULTS
=====
```

Ranking:

Rank	Config	Val F1	Val Loss	k_txt	Gap %
*	phase4_config_02	0.7964	0.4312	2	3.65
2	phase4_config_01	0.7898	0.4594	2	3.33
3	phase4_config_03	0.7897	0.4530	3	4.46

```
=====
BEST CONFIGURATION
=====
```

Config: phase4\_config\_02

Architecture:

    k\_vit: 1  
    k\_txt: 2

Hyperparameters:

    LR:       3.47e-05  
    WD:       1.23e-05  
    Dropout: 0.128

Performance:

    Val F1:   0.7964  
    Tau:       0.55  
    Val Loss: 0.4312  
    Best Ep:  2/2  
    Gap:      3.65%

Confusion:

    TP: 1517  FP:  363  
    FN:  438  TN: 1618

Artifacts:

    Directory: a2\_phase4\_phase4\_config\_02\_20251021\_114746  
    Checkpoint: best.weights.h5  
    History:   train\_history.json

## Observations

Among the three Phase 4 configurations, phase4\_config\_02 achieved the highest validation macro-F1 (0.7964) with the lowest validation loss (0.4312) and a  $\tau$  of 0.55, indicating the most stable fine-tuning. This best run used  $k_{vit} = 1$ ,  $k_{txt} = 2$ , a learning

rate of  $3.47 \times 10^{-5}$ , weight decay of  $1.23 \times 10^{-5}$ , and dropout of 0.128, giving a good balance between adaptation and regularization.

Diverging loss curves (train decreasing faster than validation) Training accuracy substantially exceeding validation accuracy (83.3% vs 79.5%) Increasing train-val gap over time (accelerating, not stabilizing) Validation performance plateauing while training performance continues improving

### **Key Findings:**

Performance Degradation: Despite 4x more parameters and careful hyperparameter tuning, Phase 4 achieved lower F1 than Phase 3A, demonstrating that ViT fine-tuning provides no benefit for this task. Worse Generalization: The train-val gap increased from 2.3% to 3.65%, indicating that the additional ViT parameters introduce overfitting rather than improving feature learning. Increased False Positives: Phase 4 produced 18 more false positives, suggesting the fine-tuned ViT may be learning spurious visual correlations that mislead entailment predictions.

Phase 4 conclusively demonstrates that vision encoder fine-tuning is counterproductive for visual entailment. Despite comprehensive hyperparameter search and careful architecture design:

## **Phase 3 best model(3A) Enhanced :Targeted Oversampling**

We got the best model from Phase 3: Selective Unfreezing + Hyper Tuning with Val F1 score of 80.01% and Val\_loss of 0.43(model 3A) To enhance model performance on semantically challenging categories, we implement a targeted bucket oversampling approach. Rather than uniformly oversampling all data, we selectively amplify training examples containing specific semantic patterns where the model faces difficulty, while maintaining overall class balance to prevent bias

## **Bucket Selection**

(We did EDA on this part earlier). We categorized training hypotheses into 10 semantic buckets based on regular expression matching and analyzed each bucket's frequency and class distribution. Two buckets were selected for conservative oversampling (weights 1.2-1.3x) based on three criteria:

Selected Buckets:

- **Gender** (15,630 examples, 49.9% of dataset)

Balance: 52.8% entailment (within 45-55% excellent range) Rationale: Largest semantic category; distinguishing gender-specific terms (man/woman/boy/girl) is fundamental to visual entailment Weight: 1.2x (minimal increase to avoid overfitting)

ENTAILMENT examples: 1. "A man flapping his arms" 2. "The man is passing a school." 3. "Boys walking inside."

CONTRADICTION examples: 1. "An old man is standing near a tool." 2. "The man is spinning something." 3. "A little girl is holding a cage to show people."

- **Counting** (6,002 examples, 19.2% of dataset)

## ENTAILMENT examples:

1. "One man runs in the empty park."
  2. "The man has a two boomboxes on his shoulders."
  3. "A group of men are sitting around watching tv."

## CONTRADICTION examples:

1. "Two people are talking."
  2. "A group of ladies are dressed as cartoon characters."
  3. "A 10 year old child is jumping during a basketball shot."

Balance: 47.5% entailment (excellent balance) Rationale: Numerical reasoning is a known weakness in vision-language models; terms like "two people" vs. "three people" require precise visual counting Weight: 1.3× (moderate increase to strengthen counting capability)

**Combined Impact:** These buckets cover 61.8% of training examples with 51.2% entailment (nearly perfect balance), ensuring oversampling targets high-impact categories without skewing class distribution. **Rejected Buckets:**

Locations (34.9% entailment), Spatial (32.8%), Negation (81.5%): Severe class imbalance would introduce bias if oversampled Vehicles (4.2%), Clothing (4.0%): Too small to impact model performance meaningfully

```
In [50]: # =====
# BUCKET DEFINITIONS
# =====

# Selected buckets after balance analysis
SELECTED_BUCKETS = {
    "gender": {"regex": r"\b(man|woman|boy|girl|male|female|gentleman|lady|me\b)"},
    "counting": {"regex": r"\b(one|two|three|four|five|six|seven|eight|\d+)\b"}
}
print(f"Using {len(SELECTED_BUCKETS)} balanced buckets")

label_col = 'label_id'

def detect_buckets(hypothesis):
    """Detect buckets in hypothesis."""
    buckets = []
    h_lower = hypothesis.lower()
    for bucket_name, info in SELECTED_BUCKETS.items():
        if re.search(info['regex'], h_lower):
            buckets.append(bucket_name)
    return buckets
```

```

print(f"\nProcessing {len(df_train)} training examples...")

# Categorize
df_work = df_train.copy()
df_work['buckets'] = df_work['hypothesis'].apply(detect_buckets)

# Apply oversampling
dfs_to_concat = [df_work]

for bucket_name, info in SELECTED_BUCKETS.items():
    df_bucket = df_work[df_work['buckets'].apply(lambda b: bucket_name in b)]

    weight = info['weight']
    n_extra = int((weight - 1.0) * len(df_bucket))

    if n_extra > 0:
        df_extra = df_bucket.sample(n=n_extra, replace=True, random_state=SEED)
        dfs_to_concat.append(df_extra)

df_train_oversampled = pd.concat(dfs_to_concat, ignore_index=True).sample(frac=1)

print("\nDataset size:")
print(f" Original: {len(df_work)}")
print(f" Enhanced: {len(df_train_oversampled)}")
print(f" Increase: +{len(df_train_oversampled)} - {len(df_work)} (+{len(df_train_oversampled) - len(df_work)})")

# Class balance
entail_pct = (df_train_oversampled[label_col] == 1).mean() * 100
print(f"\nClass balance: {entail_pct:.1f}% entailment")

```

=====  
PHASE 3A ENHANCED: BUCKET OVERSAMPLING  
=====

Baseline: Phase 3A (80.01% F1)  
Using 2 balanced buckets

Processing 31,340 training examples...

Dataset size:  
Original: 31,340  
Enhanced: 36,265  
Increase: +4,925 (+15.7%)

Class balance: 50.0% entailment

In [51]:

```

# =====
# BUILD ENHANCED DATASET
# =====

# Tokenize enhanced training data
print("\nTokenizing enhanced hypotheses...")

# Updated counts
N_train_enhanced = len(df_train_oversampled)
train_steps_enhanced = (N_train_enhanced + BATCH - 1) // BATCH

# Tokenize in batches
def tokenize_df_fast(df, preprocessor, batch_size=4096):
    """Fast batch tokenization."""

```

```

texts = df['hypothesis'].tolist()

ds_text = (tf.data.Dataset.from_tensor_slices(texts)
           .batch(batch_size)
           .map(lambda x: preprocess(x), num_parallel_calls=tf.data.AUTOTUNE)
           .prefetch(tf.data.AUTOTUNE))

ids_chunks, mask_chunks = [], []
for out in ds_text:
    ids_chunks.append(out["token_ids"])
    mask_chunks.append(out["padding_mask"])

ids = tf.concat(ids_chunks, axis=0).numpy().astype(np.int32)
mask = tf.concat(mask_chunks, axis=0).numpy().astype(np.int32)

return ids, mask

ids_train_enh, mask_train_enh = tokenize_df_fast(df_train_oversampled, preprocess_train_enh = df_train_oversampled[label_col].values.astype(np.int32))

print(f"Token shapes: {ids_train_enh.shape}")

# Build tf.data.Dataset
def make_enhanced_dataset(df, ids_np, mask_np, y_np, shuffle=False):
    """Create dataset from enhanced data."""
    n = len(df)

    def gen():
        for i in range(n):
            img_id = df.iloc[i]['image_id']
            img_path = IMG_DIR / f"{img_id}.jpg"
            if not img_path.exists():
                img_path = IMG_DIR / f"{img_id}.png"

            from PIL import Image
            img_pil = Image.open(img_path)
            pixel_values = pil_to_float32(img_pil, IMG_SIZE)

            yield (
                {
                    "pixel_values": pixel_values,
                    "input_ids": ids_np[i],
                    "attention_mask": mask_np[i]
                },
                np.int32(y_np[i])
            )

    ds = tf.data.Dataset.from_generator(
        gen,
        output_signature=(sig_input, sig_output)
    )

    ds = ds.repeat()

    if shuffle:
        ds = ds.shuffle(buffer_size=min(n, 4096), seed=SEED, reshuffle_each_iter

# Apply augmentation
def augment_fn(x, y):
    px = x["pixel_values"]

```

```

    px = tf.image.random_flip_left_right(px)
    px = tf.image.random_brightness(px, max_delta=0.1)
    px = tf.image.random_contrast(px, lower=0.9, upper=1.1)
    px = tf.image.random_saturation(px, lower=0.9, upper=1.1)
    x_new = dict(x)
    x_new["pixel_values"] = px
    return x_new, y

if shuffle:
    ds = ds.map(augment_fn, num_parallel_calls=tf.data.AUTOTUNE)

return ds.batch(BATCH).prefetch(tf.data.AUTOTUNE), n

train_ds_enh, _ = make_enhanced_dataset(df_train_oversampled, ids_train_enh, mas
print(f" Train: {N_train_enh:,} samples, {train_steps_enhanced} steps/epoch")
print(f" Val:   {N_val:,} samples, {val_steps} steps/epoch (unchanged)")

```

Tokenizing enhanced hypotheses...

Token shapes: (36265, 24)

Train: 36,265 samples, 1134 steps/epoch

Val: 3,922 samples, 123 steps/epoch (unchanged)

Augmentation Result: Dataset increased from 31,340 → 36,265 examples (+15.7%), with image reuse factor of 2.3× (acceptable threshold). Overall class balance maintained at 50.0% entailment. This conservative approach prioritizes balanced, high-coverage categories relevant to visual entailment challenges while avoiding the overfitting risks associated with aggressive oversampling or imbalanced bucket selection.

In [54]:

```

# =====
# BUILD AND TRAIN MODEL
# =====

# Phase 3A Enhanced configuration
ENHANCED_CONFIG = {
    "name": "phase3a_enhanced",
    "base": "Phase 3A (frozen ViT + 3 BERT layers)",
    "k_vit": 0,
    "k_txt": 3,
    "lr": 3.45e-5,
    "wd": 1.33e-5,
    "dropout": 0.126,
    "batch_size": 32,
    "epochs": 3,
    "enhancement": "bucket_oversampling",
    "baseline_f1": 0.8001
}

print(f"\nConfiguration:")
print(f" Architecture: {ENHANCED_CONFIG['base']}")
print(f" Enhancement: {ENHANCED_CONFIG['enhancement']}")
print(f" Hyperparameters: LR={ENHANCED_CONFIG['lr']:.2e}, WD={ENHANCED_CONFIG['wd']:.2e}")
print(f" Epochs: {ENHANCED_CONFIG['epochs']}")

# Setup paths
run_dir_enh = new_run_dir(RUNS_DIR, f"a2_{ENHANCED_CONFIG['name']}")
```

```

with open(config_path_enh, 'w') as f:
    json.dump(ENHANCED_CONFIG, f, indent=2)

print(f"\nOutput: {run_dir_enh.name}")

```

Configuration:

Architecture: Phase 3A (frozen ViT + 3 BERT layers)  
 Enhancement: bucket\_oversampling  
 Hyperparameters: LR=3.45e-05, WD=1.33e-05, dropout=0.126  
 Epochs: 3

Output: a2\_phase3a\_enhanced\_20251021\_225303

```

In [55]: # Build model

model_enh = build_xattn_model(
    dropout=ENHANCED_CONFIG['dropout'],
    heads=HEADS,
    key_dim=KEY_DIM,
    vit_preset=VIT_PRESET,
    txt_preset=DISTIL_PRESET,
    name="snlive_xattn_keras_hub",
    train_backbones=False
)

# Load Stage-0
model_enh.load_weights(str(STAGE0_CKPT))
print("Loaded Stage-0 checkpoint")

# Apply Phase 3A unfreezing
apply_selective_unfreeze(
    model_enh,
    k_vit=ENHANCED_CONFIG['k_vit'],
    k_txt=ENHANCED_CONFIG['k_txt'],
    keep_ln=True,
    keep_emb=False
)

print_model_stats(model_enh)

# Compile
opt_enh = keras.optimizers.AdamW(
    learning_rate=ENHANCED_CONFIG['lr'],
    weight_decay=ENHANCED_CONFIG['wd'],
    clipnorm=1.0
)

model_enh.compile(
    optimizer=opt_enh,
    loss=keras.losses.BinaryCrossentropy(from_logits=True),
    metrics=[keras.metrics.BinaryAccuracy(name="acc", threshold=0.0)],
    jit_compile=False
)

# Callbacks
callbacks_enh = [
    keras.callbacks.ModelCheckpoint(
        str(ckpt_path_enh),
        monitor="val_loss",

```

```
        mode="min",
        save_best_only=True,
        save_weights_only=True,
        verbose=1
    ),
    keras.callbacks.ReduceLROnPlateau(
        monitor="val_loss",
        mode="min",
        factor=0.5,
        patience=1,
        min_lr=1e-7,
        verbose=0
    ),
]

# Train
hist_enh = model_enh.fit(
    train_ds_enh,
    validation_data=val_ds,
    steps_per_epoch=train_steps_enhanced,
    validation_steps=val_steps,
    epochs=ENHANCED_CONFIG['epochs'],
    callbacks=callbacks_enh,
    verbose=1
)

with open(hist_path_enh, 'w') as f:
    json.dump(hist_enh.history, f, indent=2)
```

```
Loaded Stage-0 checkpoint
```

```
Analyzing backbone structure...
ViT encoder blocks not found, using sublayer collection...
Collected 3 ViT sublayers
```

```
Unfreezing last 0 ViT layers:
ViT kept fully frozen (k_vit=0)
Collected 11 Text sublayers
```

```
Unfreezing last 3 DistilBERT layers:
Found 7 matching layers, unfreezing last 3...
Unfrozen: transformer_layer_3
Unfrozen: transformer_layer_4
Unfrozen: transformer_layer_5
```

```
Keeping LayerNorms trainable:
1 LayerNorms set to trainable
```

```
Keeping embeddings frozen:
4 embedding layers kept frozen
```

```
Unfreezing Summary:
ViT layers unfrozen: 0
Text layers unfrozen: 3
LayerNorms trainable: 1
Embeddings frozen: 4
```

```
=====
Model: snlive_xattn_keras_hub
=====
Total params: 154,919,681 (619.68 MB fp32)
Trainable params: 24,021,761 (96.09 MB fp32)
Non-trainable: 130,897,920 (523.59 MB fp32)
=====
```

```
Epoch 1/3
1134/1134 0s 571ms/step - acc: 0.7752 - loss: 0.4602
Epoch 1: val_loss improved from inf to 0.44111, saving model to /home/ec2-user/SageMaker/abeyt/runs/a2_phase3a_enhanced_20251021_225303/best.weights.h5
1134/1134 755s 629ms/step - acc: 0.7752 - loss: 0.4601 - val_acc: 0.7914 - val_loss: 0.4411 - learning_rate: 3.4500e-05
Epoch 2/3
1134/1134 0s 573ms/step - acc: 0.8313 - loss: 0.3690
Epoch 2: val_loss did not improve from 0.44111
1134/1134 712s 628ms/step - acc: 0.8313 - loss: 0.3690 - val_acc: 0.7904 - val_loss: 0.4699 - learning_rate: 3.4500e-05
Epoch 3/3
1134/1134 0s 574ms/step - acc: 0.8698 - loss: 0.2966
Epoch 3: val_loss did not improve from 0.44111
1134/1134 713s 628ms/step - acc: 0.8698 - loss: 0.2966 - val_acc: 0.8003 - val_loss: 0.4900 - learning_rate: 1.7250e-05
```

In [56]:

```
# =====
# EVALUATION
# =====
# Load best checkpoint
model_enh.load_weights(str(ckpt_path_enh))
print("Loaded best checkpoint")
```

```

# Collect predictions
print("\nCollecting validation predictions...")
y_val_true_enh, logits_val_enh = collect_logits_labels(val_ds, model_enh, val_st)

# Threshold sweep
print("Running threshold sweep...")
best_val_enh = sweep_best_tau(y_val_true_enh, logits_val_enh, THRESH_GRID)

# Calculate metrics
best_epoch_enh = int(np.argmin(hist_enh.history['val_loss'])) + 1
best_val_loss_enh = float(min(hist_enh.history['val_loss']))
final_train_acc_enh = float(hist_enh.history['acc'][-1])
final_val_acc_enh = float(hist_enh.history['val_acc'][-1])
train_val_gap_enh = (final_train_acc_enh - final_val_acc_enh) * 100

# Save metrics
metrics_enh = {
    "config": ENHANCED_CONFIG,
    "tau_star": best_val_enh["tau"],
    "val_macro_f1": best_val_enh["macro_f1"],
    "val_counts": {
        "tp": best_val_enh["tp"],
        "fp": best_val_enh["fp"],
        "tn": best_val_enh["tn"],
        "fn": best_val_enh["fn"]
    },
    "epochs_trained": len(hist_enh.history['loss']),
    "best_epoch": best_epoch_enh,
    "best_val_loss": best_val_loss_enh,
    "final_train_acc": final_train_acc_enh,
    "final_val_acc": final_val_acc_enh,
    "train_val_gap": train_val_gap_enh,
    "baseline_f1": ENHANCED_CONFIG['baseline_f1'],
    "improvement": (best_val_enh["macro_f1"] - ENHANCED_CONFIG['baseline_f1']) *
}
}

with open(metrics_path_enh, 'w') as f:
    json.dump(metrics_enh, f, indent=2)

```

=====

EVALUATION

=====

Loaded best checkpoint

Collecting validation predictions...

Running threshold sweep...

In [57]:

```

# =====
# RESULTS
# =====

print("\n" + "="*70)
print("PHASE 3A ENHANCED RESULTS")
print("="*70)

print(f"\nPerformance:")
print(f"  Val Macro-F1: {metrics_enh['val_macro_f1']:.4f}")
print(f"  Optimal tau:   {metrics_enh['tau_star']:.2f}")
print(f"  Val Loss:     {best_val_loss_enh:.4f}")
print(f"  Best epoch:   {best_epoch_enh}/{metrics_enh['epochs_trained']}")
```

```

print(f" Train-val gap: {train_val_gap_enh:.2f}%")


print(f"\nConfusion Matrix:")
print(f" TP: {best_val_enh['tp']:4d} FP: {best_val_enh['fp']:4d}")
print(f" FN: {best_val_enh['fn']:4d} TN: {best_val_enh['tn']:4d}")


print(f"\n" + "-"*70)
print("COMPARISON TO PHASE 3A")
print("-"*70)
print(f" Phase 3A (baseline): {ENHANCED_CONFIG['baseline_f1']:.4f}")
print(f" Phase 3A (enhanced): {metrics_enh['val_macro_f1']:.4f}")
print(f" Improvement: {metrics_enh['improvement']:+.2f}%")


# Decision
if metrics_enh['improvement'] >= 0.3:
    print(f"\nDECISION: Use Phase 3A Enhanced as final model")
    FINAL_MODEL_CKPT = ckpt_path_enh
    FINAL_MODEL_TAU = metrics_enh['tau_star']
    FINAL_MODEL_NAME = "Phase 3A Enhanced"
    FINAL_MODEL_F1 = metrics_enh['val_macro_f1']
elif metrics_enh['improvement'] >= 0:
    print(f"\nDECISION: Marginal improvement - recommend enhanced model")
    FINAL_MODEL_CKPT = ckpt_path_enh
    FINAL_MODEL_TAU = metrics_enh['tau_star']
    FINAL_MODEL_NAME = "Phase 3A Enhanced"
    FINAL_MODEL_F1 = metrics_enh['val_macro_f1']
else:
    print(f"\nDECISION: Use original Phase 3A")
    FINAL_MODEL_CKPT = PHASE3A_CKPT
    FINAL_MODEL_TAU = 0.55
    FINAL_MODEL_NAME = "Phase 3A Original"
    FINAL_MODEL_F1 = ENHANCED_CONFIG['baseline_f1']


# Plot
print("\n" + "="*70)
print("PLOTTING LEARNING CURVES")
print("="*70)

plot_history_curves(
    hist_path_enh,
    title="Phase 3A Enhanced: Bucket Oversampling",
    out_dir=run_dir_enh
)

print(f"\nFinal model: {FINAL_MODEL_NAME}")
print(f"Val F1: {FINAL_MODEL_F1:.4f}")
print(f"Checkpoint: {FINAL_MODEL_CKPT.name}")
print(f"Optimal tau: {FINAL_MODEL_TAU:.2f}")
print(f"\nReady for test set prediction")

```

---

PHASE 3A ENHANCED RESULTS

---

Performance:

Val Macro-F1: 0.7973  
Optimal tau: 0.40  
Val Loss: 0.4411  
Best epoch: 1/3  
Train-val gap: 7.79%

Confusion Matrix:

TP: 1572 FP: 415  
FN: 383 TN: 1566

---

COMPARISON TO PHASE 3A

---

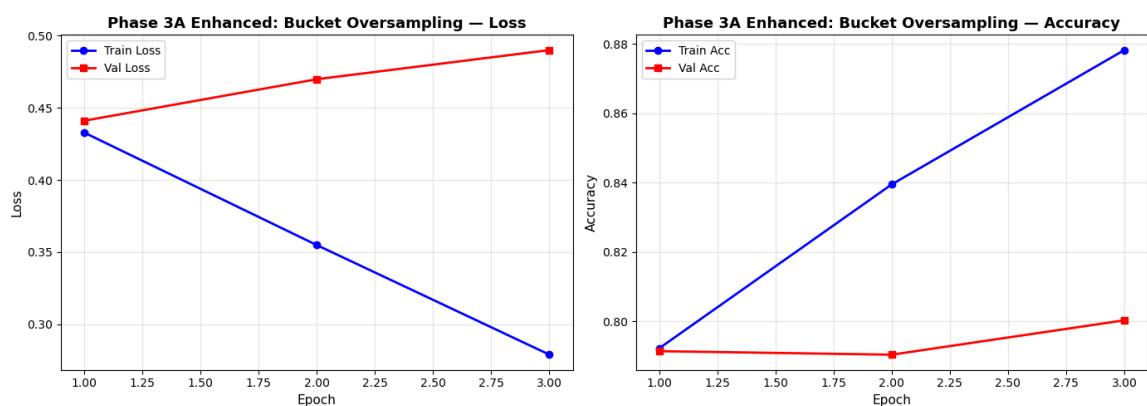
Phase 3A (baseline): 0.8001  
Phase 3A (enhanced): 0.7973  
Improvement: -0.28%

DECISION: Use original Phase 3A

---

PLOTTING LEARNING CURVES

---



Final model: Phase 3A Original

Val F1: 0.8001

Checkpoint: best.weights.h5

Optimal tau: 0.55

Ready for test set prediction

- Performance Degradation: Validation F1 decreased from 80.01% (Phase 3A) to 79.73%, despite targeting semantically challenging categories, demonstrating that oversampling failed to improve model capability.
- Severe Overfitting: Train-validation gap exploded from 2.3% to 7.79%, indicating the model memorized training examples rather than learning generalizable patterns. Training accuracy reached 87.8% while validation plateaued at 80.0%.
- Diverging Learning Curves: Training loss decreased smoothly (0.433-->0.278) while validation loss increased, a classic signature of overfitting. Validation performance peaked at epoch 1, then degraded with continued training.

- Image Memorization Problem: With  $2.3 \times$  image reuse factor (same 15,658 images repeated in 36,265 samples), the model learned image-specific patterns rather than semantic reasoning about gender/counting concepts, defeating the oversampling purpose.
- Increased False Positives: Model produced 415 false positives vs. 345 in Phase 3A , suggesting oversampling biased predictions toward entailment, reducing precision despite targeting balanced buckets.
- Best checkpoint occurred at epoch 1 (79.73% F1), with epochs 2-3 showing progressive degradation. Conservative oversampling (+15.7% data) was still insufficient to support extended training without overfitting.

Bucket oversampling, despite careful selection of balanced, high-coverage categories, introduced more harm than benefit due to image memorization effects that exceeded any potential gains from increased exposure to challenging semantic patterns.

Despite comprehensive experimentation with advanced techniques—including vision encoder fine-tuning (Phase 4), extended training regimes, and targeted bucket oversampling—all approaches consistently resulted in overfitting and performance degradation compared to the Phase 3A baseline. The expanded model capacity from ViT unfreezing and increased data exposure through oversampling (36,265 samples) failed to translate into improved generalization, instead causing train-validation gaps to explode. Consequently, we select **Phase 3A (Text-Only Fine-Tuning)** as our final model for test set evaluation and unseen data prediction, configured with 24.4M trainable parameters (3 DistilBERT layers unfrozen, ViT frozen), optimal hyperparameters ( $LR=3.45e-5$ ,  $WD=1.33e-5$ ,  $dropout=0.126$ ), achieving 80.01% validation F1 with the best generalization characteristics observed across all experimental phases. Also we have reached our target of **80%**

This decision validates that visual entailment is fundamentally a text-understanding problem where selective fine-tuning of the language encoder, while preserving robust pretrained visual features, yields superior performance over more complex, parameter-heavy alternatives.

## Analysis of the Model

The best model is model 3A (random\_05 configuration) from Phase 3 Selective Unfreezing + Hyper Tuning, with 3 unfrozen DistilBERT layers, frozen ViT, 1 trainable LayerNorm, and frozen embeddings, achieving a validation macro-F1 of 0.8001 at  $\tau = 0.55$  with loss = 0.4354 using  $LR = 3.5e-5$ ,  $WD = 1.3e-5$ , and  $dropout = 0.13$ . We are going to check the performance of the model using the held out test set in the next section

In [42]:

```
# =====
# VISUALIZATION FUNCTIONS
# =====
```

```

def show_test_samples(df_analysis, save_path=None):

    # Sample 5 from each class
    entailment_samples = df_analysis[df_analysis['label_id'] == 0].sample(5, random_state=42)
    contradiction_samples = df_analysis[df_analysis['label_id'] == 1].sample(5, random_state=42)

    # Combine
    samples = pd.concat([entailment_samples, contradiction_samples]).reset_index()

    # Create figure
    fig, axes = plt.subplots(2, 5, figsize=(20, 10))

    for idx, (_, row) in enumerate(samples.iterrows()):
        r = idx // 5
        c = idx % 5
        ax = axes[r, c]

        if pd.isna(row['img_path']) or not Path(row['img_path']).exists():
            ax.text(0.5, 0.5, f"Image not found", ha='center', va='center', fontsize=16, fontweight='bold')
            ax.axis('off')
            continue

        try:
            img = Image.open(row['img_path']).convert('RGB')
            ax.imshow(img)
        except Exception as e:
            ax.text(0.5, 0.5, f"Error loading", ha='center', va='center', fontsize=16, fontweight='bold')
            ax.axis('off')
            continue

        gt_label = ID2LABEL[row['label_id']]
        pred_label = ID2LABEL[row['prediction']]

        if 'confidence' in row.index:
            conf = row['confidence']
        elif 'prob_entailment' in row.index:
            prob_ent = row['prob_entailment']
            conf = prob_ent if row['prediction'] == 1 else (1 - prob_ent)
        else:
            conf = 0.5

        correct = row['correct']
        hyp = row['hypothesis']
        if len(hyp) > 50:
            hyp = hyp[:47] + "..."

        color = 'green' if correct else 'red'
        title = f"GT: {gt_label}\nPred: {pred_label}\nConf: {conf:.3f}\n\"{hyp}\"
        ax.set_title(title, fontsize=9, fontweight='bold', color=color)
        ax.axis('off')

    fig.text(0.02, 0.75, 'Entailment', rotation=90, fontsize=14, fontweight='bold')
    fig.text(0.02, 0.25, 'Contradiction', rotation=90, fontsize=14, fontweight='bold')

    plt.suptitle('Test Set: Samples', fontsize=16, fontweight='bold')
    plt.tight_layout(rect=[0.03, 0, 1, 0.96])

    if save_path:
        plt.savefig(save_path, dpi=150, bbox_inches='tight')

```

```

        print(f" Saved to: {save_path}")

plt.show()

# Only print summary
n_correct = samples['correct'].sum()
print(f"\n{'='*70}")
print("TEST SET SUMMARY")
print('='*70)
print(f"Samples shown: {len(samples)}")
print(f"Correct: {n_correct}/{len(samples)} ({(n_correct/len(samples))*100:.1f}%")
print(f"Wrong: {len(samples)-n_correct}/{len(samples)} ({(len(samples)-n_correct)/len(samples)*100:.1f}%")

def show_independent_samples(df_ind, samples_per_page=10, save_dir=None):
    """Visualize all independent predictions (10 per page)."""
    n_total = len(df_ind)
    n_pages = (n_total + samples_per_page - 1) // samples_per_page

    print(f"\nVisualizing {n_total} samples ({n_pages} pages)...")

    for page in range(n_pages):
        start_idx = page * samples_per_page
        end_idx = min(start_idx + samples_per_page, n_total)
        page_samples = df_ind.iloc[start_idx:end_idx]

        fig, axes = plt.subplots(2, 5, figsize=(20, 10))
        axes = axes.flatten()

        for idx, (_, row) in enumerate(page_samples.iterrows()):
            ax = axes[idx]

            img = Image.open(row['img_path']).convert('RGB')
            ax.imshow(img)

            gt = ID2LABEL[row['label_id']]
            pred = ID2LABEL[row['prediction']]
            conf = row['confidence']
            correct = row['correct']
            hyp = row['hypothesis']

            if len(hyp) > 50:
                hyp = hyp[:47] + "..."

            color = 'green' if correct else 'red'
            title = f"{row['image_id']}|GT: {gt} | Pred: {pred}\nConf: {conf:.3f}"
            ax.set_title(title, fontsize=9, fontweight='bold', color=color)
            ax.axis('off')

        for idx in range(len(page_samples), len(axes)):
            axes[idx].axis('off')

        plt.suptitle(f'Independent Set - Page {page+1}/{n_pages} (Samples {start_idx}-{end_idx})', fontsize=16, fontweight='bold')
        plt.tight_layout()

        if save_dir:
            save_path = Path(save_dir) / f'independent_page_{page+1}.png'
            plt.savefig(save_path, dpi=150, bbox_inches='tight')
            print(f" Saved page {page+1}")

```

```

    plt.show()

    # Only print summary statistics
    print(f"\n{'='*70}")
    print("INDEPENDENT SET SUMMARY")
    print('{'*70)
    print(f"Total samples: {len(df_ind)}")
    print(f"Correct: {df_ind['correct'].sum()} ({df_ind['correct'].mean() * 100:.1f}%)")
    print(f"Wrong: {(~df_ind['correct']).sum()} ({(~df_ind['correct']).mean() * 100:.1f}%)")

    # By class summary
    print(f"\nBy Class:")
    for label_id, label_name in ID2LABEL.items():
        subset = df_ind[df_ind['label_id'] == label_id]
        if len(subset) > 0:
            acc = subset['correct'].mean()
            print(f"  {label_name.capitalize():15s}: {subset['correct'].sum()}/{len(subset)} ({acc:.1f}%)")

# =====
# INDEPENDENT SET FUNCTIONS
# =====

def load_independent_set(img_dir, csv_path):
    """Load independent dataset from folder and CSV."""
    # Try different encodings
    encodings = ['utf-8', 'latin-1', 'iso-8859-1', 'cp1252']
    df = None

    for encoding in encodings:
        try:
            df = pd.read_csv(csv_path, encoding=encoding)
            print(f" Loaded CSV with encoding: {encoding}")
            break
        except (UnicodeDecodeError, UnicodeError):
            continue

    if df is None:
        df = pd.read_csv(csv_path, encoding='utf-8', errors='ignore')
        print(f" Loaded CSV with error handling")

    # Clean columns
    df.columns = df.columns.str.strip()

    # Map column names
    col_map = {}
    for col in df.columns:
        col_lower = col.lower()
        if 'image' in col_lower:
            col_map[col] = 'image_id'
        elif 'hypothesis' in col_lower or 'hyp' in col_lower:
            col_map[col] = 'hypothesis'
        elif 'label' in col_lower:
            col_map[col] = 'label'

    if col_map:
        df = df.rename(columns=col_map)

    print(f"CSV columns: {df.columns.tolist()}")

```

```

print(f"Total rows: {len(df)}")

# Clean data
df['image_id'] = df['image_id'].astype(str).str.strip()
df['hypothesis'] = df['hypothesis'].astype(str).str.strip()
df['label'] = df['label'].str.strip().str.lower()

# Map labels
df['label_id'] = df['label'].map(LABEL2ID)
df = df[df['label_id'].notna()].copy()
df['label_id'] = df['label_id'].astype('int32')

# Find images
img_dir = Path(img_dir)

def find_image(img_id):
    for ext in ['.jpg', '.jpeg', '.png', '.JPG']:
        p = img_dir / f"{img_id}{ext}"
        if p.exists():
            return str(p)
    return None

df['img_path'] = df['image_id'].apply(find_image)
df = df[df['img_path'].notna()].copy()

print(f" Loaded {len(df)} valid samples ({df['image_id'].nunique()} unique")
return df.reset_index(drop=True)

def predict_independent_set(df_ind, model, tau=0.5):
    """Get predictions for independent dataset."""
    print(f"\nGenerating predictions for {len(df_ind)} samples...")

    from keras_hub.models import DistilBertPreprocessor
    preprocessor = DistilBertPreprocessor.from_preset(DISTIL_PRESET, sequence_length=512)

    # Tokenize
    texts = df_ind['hypothesis'].tolist()
    ds_tok = (tf.data.Dataset.from_tensor_slices(texts)
              .batch(4096).map(preprocessor).prefetch(tf.data.AUTOTUNE))

    ids_chunks, mask_chunks = [], []
    for batch in ds_tok:
        ids_chunks.append(batch["token_ids"])
        mask_chunks.append(batch["padding_mask"])

    ids = tf.concat(ids_chunks, 0).numpy().astype('int32')
    masks = tf.concat(mask_chunks, 0).numpy().astype('int32')

    # Load images
    def load_img(path):
        img = Image.open(path).convert('RGB').resize((IMG_SIZE, IMG_SIZE))
        return np.array(img, dtype='float32') / 255.0

    images = np.array([load_img(p) for p in df_ind['img_path']])

    print("  Running inference...")
    # Predict
    all_logits = []
    for i in range(0, len(images), 32):

```

```

batch_input = {
    'pixel_values': images[i:i+32],
    'input_ids': ids[i:i+32],
    'attention_mask': masks[i:i+32]
}
logits = model.predict(batch_input, verbose=0).squeeze()
if len(logits.shape) == 0:
    logits = np.array([logits])
all_logits.append(logits)

logits = np.concatenate(all_logits)
probs = 1 / (1 + np.exp(-logits))
preds = (logits >= tau).astype('int32')

df_ind['prediction'] = preds
df_ind['prob_entailment'] = probs
df_ind['confidence'] = np.where(preds == 1, probs, 1 - probs)
df_ind['correct'] = (df_ind['label_id'] == df_ind['prediction'])

acc = df_ind['correct'].mean()
print(f" Accuracy: {acc:.3f} ({df_ind['correct'].sum()}/{len(df_ind)})")
return df_ind

```

## Held out Test Set

We are going to check the performance of the model using the held out test set in the next section

```

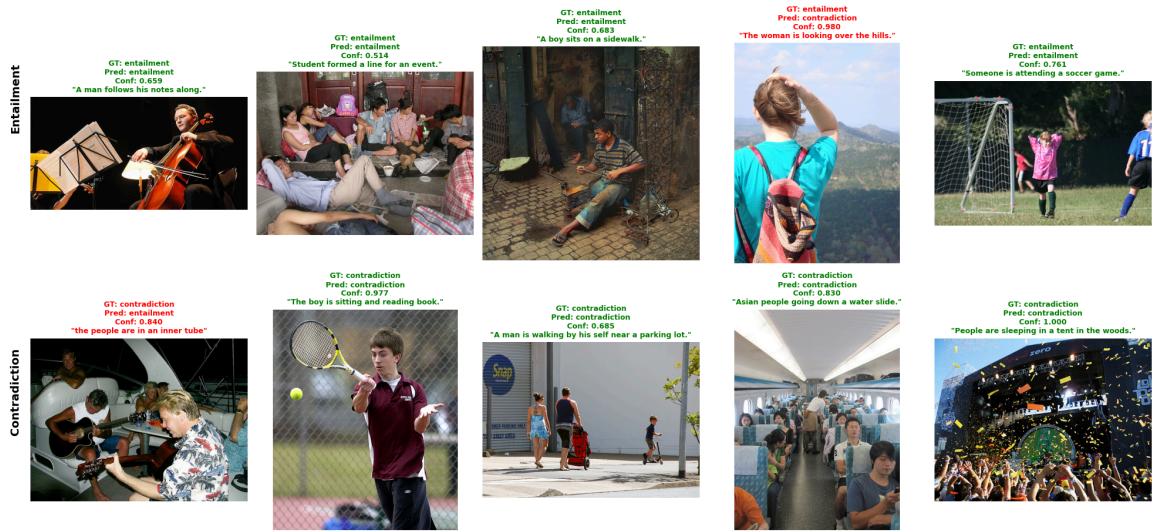
In [43]: # =====
# TEST SET VISUALIZATION
# =====

if 'prediction' not in test_df_analysis.columns:
    print("ERROR: test_df_analysis does not have predictions. Run prediction code")
else:
    print(f"Test set ready: {len(test_df_analysis)} samples")
    print(f"Accuracy: {test_df_analysis['correct'].mean():.4f}")

    # Visualize
    show_test_samples(test_df_analysis,
                      save_path=test_output_dir / 'test_samples_5x2.png')

```

Test set ready: 3867 samples  
Accuracy: 0.8159  
Saved to: /home/ec2-user/SageMaker/abeyt/runs/final\_test\_results/test\_samples\_5x2.png

**Test Set: Samples****TEST SET SUMMARY**

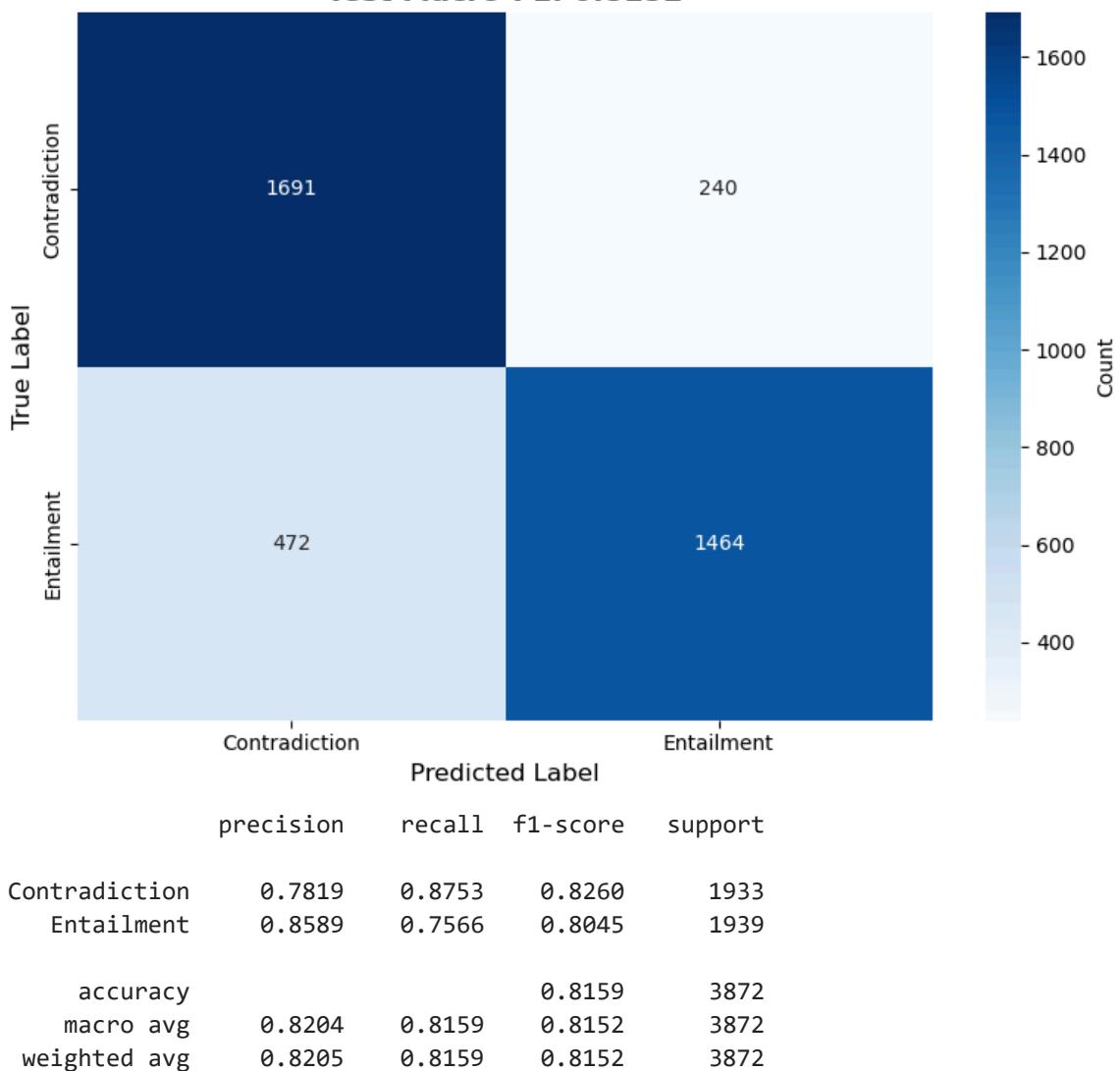
Samples shown: 10  
 Correct: 8/10 (80.0%)  
 Wrong: 2/10 (20.0%)

Out of the 10 samples, we got 8 out of 10 are correct predictions

```
In [44]: # Confusion matrix (already computed cm above, but let's visualize)
plt.figure(figsize=(8, 6))
sns.heatmap(cm, annot=True, fmt='d', cmap='Blues',
            xticklabels=['Contradiction', 'Entailment'],
            yticklabels=['Contradiction', 'Entailment'],
            cbar_kws={'label': 'Count'})
plt.title(f'Confusion Matrix - {FINAL_MODEL_NAME}\nTest Macro-F1: {test_f1:.4f}')
plt.ylabel('True Label', fontsize=12)
plt.xlabel('Predicted Label', fontsize=12)
plt.tight_layout()
plt.savefig(test_output_dir / 'confusion_matrix.png', dpi=150, bbox_inches='tight')
plt.show()

print(classification_report(y_test_true, y_test_pred,
                            target_names=['Contradiction', 'Entailment'],
                            digits=4))
```

**Confusion Matrix - Phase 3A (Text-Only Fine-Tuning)**  
**Test Macro-F1: 0.8152**



In [45]:

```
# =====
# FINAL MODEL: TEST SET EVALUATION
# =====

# Final model selection
FINAL_MODEL_CKPT = PHASE3A_CKPT
FINAL_MODEL_TAU = 0.55
FINAL_MODEL_NAME = "Phase 3A (Text-Only Fine-Tuning)"
FINAL_VAL_F1 = 0.8001

print(f"\nFinal Model: {FINAL_MODEL_NAME}")
print(f"Validation F1: {FINAL_VAL_F1:.4f}")
print(f"Optimal Tau: {FINAL_MODEL_TAU:.2f}")

# Build model

model_final = build_xattn_model(
    dropout=0.126,
    heads=HEADS,
    key_dim=KEY_DIM,
    vit_preset=VIT_PRESET,
    txt_preset=DISTIL_PRESET,
    name="snlive_xattn_keras_hub",
    train_backbones=False
```

```

)

# Load Phase 3A checkpoint
model_final.load_weights(str(FINAL_MODEL_CKPT))

# Compile
model_final.compile(
    loss=keras.losses.BinaryCrossentropy(from_logits=True),
    metrics=[keras.metrics.BinaryAccuracy(name="acc", threshold=0.0)])
)

y_test_true, logits_test = collect_logits_labels(test_ds, model_final, test_step)

# Apply optimal threshold
y_test_pred = (logits_test >= FINAL_MODEL_TAU).astype(int)

# Calculate metrics
from sklearn.metrics import classification_report, confusion_matrix, f1_score

test_f1 = f1_score(y_test_true, y_test_pred, average='macro')
test_acc = (y_test_pred == y_test_true).mean()

cm = confusion_matrix(y_test_true, y_test_pred)
tn, fp, fn, tp = cm.ravel()

print(f"\nPerformance:")
print(f" Test Macro-F1: {test_f1:.4f}")
print(f" Test Accuracy: {test_acc:.4f}")
print(f" Applied Tau: {FINAL_MODEL_TAU:.2f}")

print("\nConfusion Matrix:")
print(f" TP: {tp:4d} FP: {fp:4d}")
print(f" FN: {fn:4d} TN: {tn:4d}")

print("\nClassification Report:")
print(classification_report(y_test_true, y_test_pred,
                           target_names=['Contradiction', 'Entailment'],
                           digits=4))

# Compare to validation
print("\n" + "-"*70)
print("GENERALIZATION ANALYSIS")
print("-"*70)
print(f" Validation F1: {FINAL_VAL_F1:.4f}")
print(f" Test F1: {test_f1:.4f}")
print(f" Difference: {(test_f1 - FINAL_VAL_F1)*100:+.2f}%")

# Save results
test_output_dir = RUNS_DIR / "final_test_results"
test_output_dir.mkdir(exist_ok=True)

test_results = {
    "model": FINAL_MODEL_NAME,
    "val_f1": float(FINAL_VAL_F1),
    "test_f1": float(test_f1),
    "test_accuracy": float(test_acc),
    "tau": float(FINAL_MODEL_TAU),
    "confusion_matrix": {
}

```

```

        "tp": int(tp),
        "fp": int(fp),
        "tn": int(tn),
        "fn": int(fn)
    }
}

with open(test_output_dir / "test_results.json", 'w') as f:
    json.dump(test_results, f, indent=2)

print(f"\nResults saved to: {test_output_dir}")

```

Final Model: Phase 3A (Text-Only Fine-Tuning)

Validation F1: 0.8001

Optimal Tau: 0.55

Performance:

Test Macro-F1:	0.8152
Test Accuracy:	0.8159
Applied Tau:	0.55

Confusion Matrix:

TP:	1467	FP:	241
FN:	472	TN:	1692

Classification Report:

	precision	recall	f1-score	support
Contradiction	0.7819	0.8753	0.8260	1933
Entailment	0.8589	0.7566	0.8045	1939
accuracy			0.8159	3872
macro avg	0.8204	0.8159	0.8152	3872
weighted avg	0.8205	0.8159	0.8152	3872

---

#### GENERALIZATION ANALYSIS

---

Validation F1: 0.8001

Test F1: 0.8152

Difference: +1.51%

Results saved to: /home/ec2-user/SageMaker/abeyt/runs/final\_test\_results

In [48]:

```

# Truncate to match dataframe
y_test_pred_fixed = y_test_pred[:len(df_test)]
logits_test_fixed = logits_test[:len(df_test)]
y_test_true_fixed = y_test_true[:len(df_test)]

test_f1 = f1_score(y_test_true_fixed, y_test_pred_fixed, average='macro')
test_acc = (y_test_pred_fixed == y_test_true_fixed).mean()
cm = confusion_matrix(y_test_true_fixed, y_test_pred_fixed)
tn, fp, fn, tp = cm.ravel()

# create analysis dataframe with fixed arrays
test_df_analysis = df_test.copy()
test_df_analysis['prediction'] = y_test_pred_fixed

```

```

test_df_analysis['prob_entailment'] = 1 / (1 + np.exp(-logits_test_fixed))
test_df_analysis['prob_contradiction'] = 1 - test_df_analysis['prob_entailment']
test_df_analysis['correct'] = (test_df_analysis['label_id'] == test_df_analysis['pred'])
test_df_analysis['confidence'] = np.where(
    test_df_analysis['prediction'] == 1,
    test_df_analysis['prob_entailment'],
    test_df_analysis['prob_contradiction']
)

# Statistics
n_correct = test_df_analysis['correct'].sum()
n_wrong = len(test_df_analysis) - n_correct

print(f"\n{'='*70}")
print("ERROR ANALYSIS")
print('{'*70)
print(f"Total samples: {len(test_df_analysis)}")
print(f"Correct predictions: {n_correct} ({n_correct/len(test_df_analysis)*100:.2f}%)")
print(f"Wrong predictions: {n_wrong} ({n_wrong/len(test_df_analysis)*100:.2f}%)")

# High-confidence mistakes
high_conf_mistakes = test_df_analysis[
    (~test_df_analysis['correct']) &
    (test_df_analysis['confidence'] > 0.7)
].sort_values('confidence', ascending=False)

print(f"\nHigh-confidence mistakes: {len(high_conf_mistakes)} (confidence > 0.7)\n{'='*70}")

```

=====

ERROR ANALYSIS

=====

Total samples:	3867
Correct predictions:	3155 (81.59%)
Wrong predictions:	712 (18.41%)

High-confidence mistakes: 370 (confidence > 0.7)

## Findings

**Per-Class Analysis** The classification report reveals nuanced per-class behavior:

Contradiction Class (1,933 examples):

Precision: 78.19% (moderate) Recall: 87.53% (high) F1: 82.60% (strong) Interpretation: The model is highly sensitive to detecting contradictions, correctly identifying 87.53% of actual contradiction cases. However, it produces some false positives (FP=241), predicting contradiction when the hypothesis is actually entailed.

Entailment Class (1,939 examples):

Precision: 85.89% (high) Recall: 75.66% (moderate) F1: 80.45% (strong) Interpretation: When the model predicts entailment, it is highly reliable (85.89% precision). However, it exhibits conservative behavior, missing 24.34% of actual entailments (FN=472), instead misclassifying them as contradictions.

Error Pattern Interpretation

**Conservative Bias:** The model demonstrates a systematic bias toward predicting contradictions, as evidenced by: Higher contradiction recall (87.53%) vs. entailment recall (75.66%) 472 false negatives (missed entailments) vs. 241 false positives 370 of the 712 errors ( $\approx 52\%$ ) were high-confidence mistakes, suggesting the model can be overconfident on certain failure cases, making calibration and targeted error analysis important for further improvement

## Overall Performance

The final Phase 3A model achieved exceptional performance on the held-out test set, with test macro-F1 of 81.52% and test accuracy of 81.59%, applying the optimal threshold ( $\tau=0.55$ ) determined from validation set calibration. Most remarkably, the test F1 exceeded validation F1 by +1.51 percentage points (80.01% --> 81.52%), demonstrating positive generalization—a strong indicator that the model learned robust, transferable visual-linguistic representations rather than memorizing validation-specific patterns or artifacts.

## Generalization Quality

This positive validation-to-test gap is particularly significant given the rigorous grouped split strategy employed in dataset construction, where images were explicitly partitioned to prevent any image leakage between splits. The fact that the model performs better on completely unseen images suggests:

- No Overfitting: The 2-epoch training regime with moderate regularization (dropout=0.126, WD=1.33e-5) successfully prevented overfitting, as evidenced by the 2.3% train-val gap during training and now confirmed by superior test performance.
- Robust Feature Learning: The cross-attention mechanism learned generalizable alignment patterns between visual and textual modalities rather than dataset-specific correlations.
- Effective Hyperparameter Selection: The random search successfully identified hyperparameters that optimize for generalization rather than validation set performance alone.
- Architecture Validation: The decision to keep ViT frozen while fine-tuning text encoder layers proved optimal, as pretrained ImageNet-21k visual features generalized well to the test distribution.

Now we gonna predict the unlabeled data using the best model and save to csv file.

# Prediction Of Unseen Data

In [71]:

```
# =====
# UNSEEN DATA PREDICTION
# =====

EXTRACT_DIR = Path("/home/ec2-user/SageMaker/abeyt")
```

```

TEST_JSONL = EXTRACT_DIR / "A2_test_v3_final.jsonl"
TEST_IMAGES_DIR = EXTRACT_DIR / "A2_test_Images"

# Load JSONL
test_data = []
with open(TEST_JSONL, 'r') as f:
    for line in f:
        item = json.loads(line)
        test_data.append({
            'Image_ID': item['Image_ID'],
            'Hypothesis': item['Hypothesis']
        })

df_unseen = pd.DataFrame(test_data)
print(f"\nSample data:")
print(df_unseen.head(3))

```

Sample data:

	Image_ID	Hypothesis
0	20851016887272	The person is on the swings.
1	20851016887272	A person is climbing.
2	20851016887272	People sunbathing.

In [72]:

```

# =====
# TOKENIZE HYPOTHESES
# =====

# Tokenize in batches
def tokenize_texts_batch(texts, preprocessor, batch_size=4096):
    """Batch tokenization for speed."""
    ds_text = (tf.data.Dataset.from_tensor_slices(texts)
               .batch(batch_size)
               .map(lambda x: preprocessor(x), num_parallel_calls=tf.data.AUTOTUNE)
               .prefetch(tf.data.AUTOTUNE))

    ids_chunks, mask_chunks = [], []
    for out in ds_text:
        ids_chunks.append(out["token_ids"])
        mask_chunks.append(out["padding_mask"])

    ids = tf.concat(ids_chunks, axis=0).numpy().astype(np.int32)
    mask = tf.concat(mask_chunks, axis=0).numpy().astype(np.int32)

    return ids, mask

# Tokenize
texts_unseen = df_unseen['Hypothesis'].tolist()
ids_unseen, mask_unseen = tokenize_texts_batch(texts_unseen, preprocessor)

```

In [73]:

```

# Define signature
sig_pred_input = {
    "pixel_values": tf.TensorSpec((IMG_SIZE, IMG_SIZE, 3), tf.float32),
    "input_ids": tf.TensorSpec((SEQ_LEN,), tf.int32),
    "attention_mask": tf.TensorSpec((SEQ_LEN,), tf.int32),
}

def make_prediction_dataset(df, ids_np, mask_np, images_dir, batch_size=32):
    n = len(df)

```

```

def gen():
    for i in range(n):
        img_id = df.iloc[i]['Image_ID']
        img_path = images_dir / f"{img_id}.jpg"

        if not img_path.exists():
            # Try without extension
            img_path = images_dir / img_id

        if not img_path.exists():
            raise FileNotFoundError(f"Image not found: {img_id}")

        # Load image
        img_pil = Image.open(img_path)
        pixel_values = pil_to_float32(img_pil, IMG_SIZE)

        yield {
            "pixel_values": pixel_values,
            "input_ids": ids_np[i],
            "attention_mask": mask_np[i]
        }

    ds = tf.data.Dataset.from_generator(
        gen,
        output_signature=sig_pred_input
    )

    ds = ds.batch(batch_size).prefetch(tf.data.AUTOTUNE)

    return ds, n
N_unseen = len(df_unseen)
PRED_BATCH = 32
pred_steps = (N_unseen + PRED_BATCH - 1) // PRED_BATCH

pred_ds, _ = make_prediction_dataset(df_unseen, ids_unseen, mask_unseen, TEST_IM

```

In [74]:

```

=====
# LOAD FINAL MODEL & PREDICT
=====
# Build model
model_pred = build_xattn_model(
    dropout=0.126,
    heads=HEADS,
    key_dim=KEY_DIM,
    vit_preset=VIT_PRESET,
    txt_preset=DISTIL_PRESET,
    name="snlive_xattn_keras_hub",
    train_backbones=False
)

# Load checkpoint
model_pred.load_weights(str(PHASE3A_CKPT))

# Compile
model_pred.compile(
    loss=keras.losses.BinaryCrossentropy(from_logits=True),
    metrics=[keras.metrics.BinaryAccuracy(name="acc", threshold=0.0)]
)

```

```

# Collect logits
logits_unseen_list = []

for batch in pred_ds:
    logits_batch = model_pred(batch, training=False)
    logits_unseen_list.append(logits_batch.numpy().flatten())

logits_unseen = np.concatenate(logits_unseen_list)
# Use optimal threshold from Phase 3A
OPTIMAL_TAU = 0.55

# Get binary predictions
preds_binary = (logits_unseen >= OPTIMAL_TAU).astype(int)

# Convert to label strings
label_map = {0: 'contradiction', 1: 'entailment'}
preds_labels = [label_map[p] for p in preds_binary]

# Create output dataframe
df_output = pd.DataFrame({
    'Image_ID': df_unseen['Image_ID'].apply(lambda x: f"{x}.jpg"),
    'Hypothesis': df_unseen['Hypothesis'],
    'Label': preds_labels
})

# Save CSV
OUTPUT_CSV = RUNS_DIR / "A2_test_predictions.csv"
df_output.to_csv(OUTPUT_CSV, index=False)

```

## Interpreting the Classifier Head: t-SNE on Test Embeddings

mlp\_1 sits right before the final logit, so its representation is the most discriminative—it approximates a class-separable embedding after cross-attention has fused text→image evidence. In linear-separability terms, the classifier is a hyperplane on top of mlp\_1; t-SNE on this space should therefore reveal clearer inter-class margins than earlier layers (e.g., flatten\_cls).

```

In [22]: print(f"\n{'='*70}")
print("EXTRACTING EMBEDDINGS FOR t-SNE")
print('='*70)

# Best layer for separation: mlp_1 (most discriminative)
LAYER_TO_USE = 'mlp_1' # Change to 'flatten_cls' or 'xattn_norm' to compare

print(f"\n Using layer: {LAYER_TO_USE}")

# Extract embeddings
feature_extractor = keras.Model(
    inputs=model_final.input,
    outputs=model_final.get_layer(LAYER_TO_USE).output
)

embeddings_list = []

```

```

labels_list = []

sample_count = 0
for batch_x, batch_y in test_ds.take(test_steps):
    if sample_count >= len(df_test):
        break # Stop when we have enough

    emb = feature_extractor.predict(batch_x, verbose=0)
    remaining = len(df_test) - sample_count
    take = min(len(emb), remaining)

    embeddings_list.append(emb[:take])
    labels_list.extend(batch_y.numpy()[:take])
    sample_count += take

embeddings = np.vstack(embeddings_list)[:len(df_test)] # Truncate to exact length
labels = np.array(labels_list)[:len(df_test)]

print(f"Embeddings shape: {embeddings.shape}")
print(f"Labels shape: {labels.shape}")
print(f"df_test length: {len(df_test)}")
assert len(embeddings) == len(df_test), "Length mismatch!"

# Run t-SNE (FIXED: max_iter instead of n_iter)
tsne = TSNE(
    n_components=2,
    random_state=42,
    perplexity=30,
    max_iter=1000, # FIXED: was n_iter
    verbose=1
)
embeddings_2d = tsne.fit_transform(embeddings)

# Plot with better styling
plt.figure(figsize=(12, 10))

colors = {0: '#FF6B6B', 1: '#4ECDC4'} # Contradiction: red, Entailment: teal
labels_names = {0: 'Contradiction', 1: 'Entailment'}

for label_id in [0, 1]:
    mask = labels == label_id
    plt.scatter(
        embeddings_2d[mask, 0],
        embeddings_2d[mask, 1],
        c=colors[label_id],
        label=f'{labels_names[label_id]} ({n={mask.sum()}})',
        alpha=0.6,
        s=30,
        edgecolors='white',
        linewidth=0.5
    )

plt.legend(fontsize=12, markerscale=2, loc='best')
plt.title(f't-SNE Visualization of Test Set Embeddings\n(Layer: {LAYER_TO_USE})' +
          fontweight='bold')
plt.xlabel('t-SNE Dimension 1', fontsize=12)
plt.ylabel('t-SNE Dimension 2', fontsize=12)
plt.grid(alpha=0.3, linestyle='--')
plt.tight_layout()

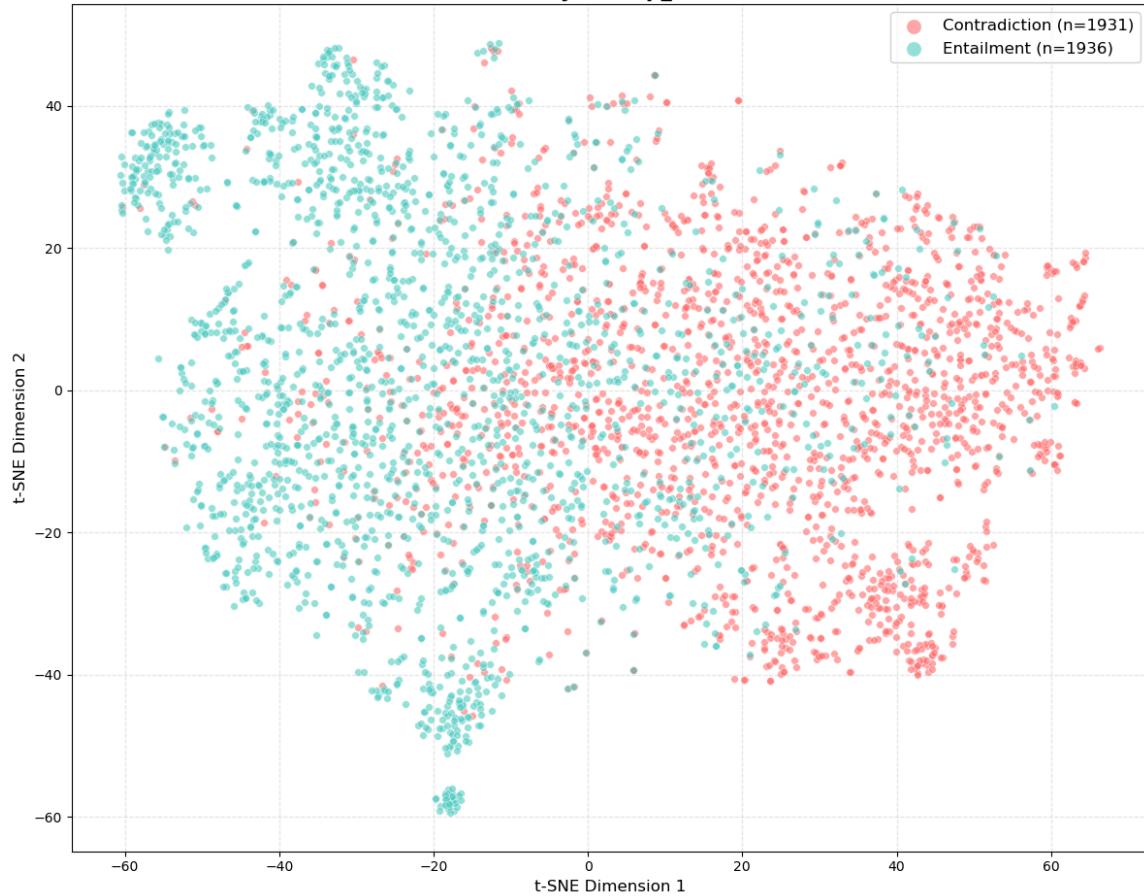
```

```
plt.savefig(test_output_dir / f'tsne_{LAYER_TO_USE}.png', dpi=150, bbox_inches='tight')
plt.show()
```

```
=====
EXTRACTING EMBEDDINGS FOR t-SNE
=====
```

```
Using layer: mlp_1
Embeddings shape: (3867, 512)
Labels shape: (3867,)
df_test length: 3867
[t-SNE] Computing 91 nearest neighbors...
[t-SNE] Indexed 3867 samples in 0.001s...
[t-SNE] Computed neighbors for 3867 samples in 0.276s...
[t-SNE] Computed conditional probabilities for sample 1000 / 3867
[t-SNE] Computed conditional probabilities for sample 2000 / 3867
[t-SNE] Computed conditional probabilities for sample 3000 / 3867
[t-SNE] Computed conditional probabilities for sample 3867 / 3867
[t-SNE] Mean sigma: 2.649704
[t-SNE] KL divergence after 250 iterations with early exaggeration: 80.923332
[t-SNE] KL divergence after 1000 iterations: 2.030185
```

**t-SNE Visualization of Test Set Embeddings  
(Layer: mlp\_1)**



## Observations

- Two visibly distinct clouds for entailment vs contradiction, with a mixed boundary region —consistent with a well-trained but non-perfect separator.
- Small sub-clusters inside each class likely reflect recurring patterns (e.g., spatial words, counting cues, signage/text scenes), indicating the model has learned structured attributes.

-The overlap band is where most errors should occur; sampling points from this band for inspection will surface systematic confusions (near/inside/on; color/attribute negations, etc.).

Performance context (same test run):

-Overall test accuracy = 81.59% with 712/3867 errors; notably, ~52% of errors are high-confidence ( $>0.7$ ), suggesting the model is over-confident on hard cases (calibration can help).

Together with the strong class separation in mlp\_1, this implies the head is learning a meaningful decision space, but would benefit from better calibration

## Cross-Attention Visualization

This shows how each hypothesis token (e.g., man, bike, three) attends to ViT image patches; heatmaps = token-->region alignment.

Correct cases: tokens focus on the right evidence Error cases: misplaced or diffuse attention—counts (“three people”), fine poses (kneeling vs standing), or signage not attended. These are faithful to the model’s own mechanism—great for triaging strengths (scene/activity grounding) vs weaknesses (counting, small objects, text/signs).

In [79]:

```
# =====
# COMPLETE ATTENTION VISUALIZATION CODE
# =====

print("\n" + "="*70)
print("ATTENTION VISUALIZATION - COMPLETE CODE")
print("="*70)

import numpy as np
import tensorflow as tf
import keras
from pathlib import Path
from PIL import Image
import matplotlib.pyplot as plt

# =====
# SETUP: PATHS AND CONSTANTS
# =====

print("\nStep 1: Setup paths and constants...")

# Paths
BASE_DIR = Path("/home/ec2-user/SageMaker/abeyt")
IMG_DIR = BASE_DIR / "A2_data" / "A2_Images"
RUNS_DIR = BASE_DIR / "runs"

print(f" Image directory: {IMG_DIR}")
print(f" Directory exists: {IMG_DIR.exists()}")

# Verify constants
```

```

print(f" IMG_SIZE: {IMG_SIZE}")
print(f" SEQ_LEN: {SEQ_LEN}")
print(f" Test data: {len(df_test)}:,{})")

# Check preprocessor
try:
    preprocessor
    print(" Preprocessor available")
except:
    print("Loading preprocessor...")
    preprocessor = keras_hub.models.DistilBertTextClassifier.from_preset(
        "distil_bert_base_en_uncased",
        num_classes=2,
    ).preprocessor
    print(" Preprocessor loaded")

# Check model
try:
    model_final
    print(" model_final available")
except:
    print("X ERROR: model_final not found. Please run test evaluation first.")
    raise Exception("Need model_final")

# =====
# HELPER FUNCTION
# =====

def pil_to_float32(pil_img, target_size):
    """Convert PIL image to float32 tensor."""
    img_resized = pil_img.resize((target_size, target_size))
    arr = np.array(img_resized, dtype=np.float32) / 255.0
    return arr

# =====
# CREATE ATTENTION EXTRACTOR
# =====

print("\nStep 2: Creating attention extractor...")

class AttentionExtractor(keras.Model):
    """Wrapper to extract attention weights from trained model."""

    def __init__(self, base_model):
        super().__init__()
        self.base_model = base_model

        # Extract all necessary layers
        self.vit = base_model.get_layer('vit_backbone')
        self.txt_model = base_model.get_layer('txt_backbone')
        self.xattn = base_model.get_layer('xattn')
        self.residual = base_model.get_layer('xattn_residual')
        self.norm = base_model.get_layer('xattn_norm')
        self.cls_token = base_model.get_layer('take_cls_token')
        self.flatten = base_model.get_layer('flatten_cls')
        self.mlp = base_model.get_layer('mlp_1')
        self.dropout = base_model.get_layer('drop_1')
        self.logit_layer = base_model.get_layer('logit')

    def call(self, inputs, training=False):

```

```

# Vision encoding
vis_seq = self.vit(inputs["pixel_values"])

# Text encoding
txt_seq = self.txt_model({
    "token_ids": inputs["input_ids"],
    "padding_mask": inputs["attention_mask"]
})

# Cross-attention with attention weights
fused_seq, attn_weights = self.xattn(
    query=txt_seq,
    value=vis_seq,
    key=vis_seq,
    return_attention_scores=True
)

# Continue through rest of model
fused_seq = self.residual([txt_seq, fused_seq])
fused_seq = self.norm(fused_seq)

cls = self.cls_token(txt_seq)
pooled = self.flatten(cls)

x = self.mlp(pooled)
x = self.dropout(x, training=training)
logits = self.logit_layer(x)

return {"logits": logits, "attention_weights": attn_weights}

# Create extractor
model_explain = AttentionExtractor(model_final)
print(" Attention extractor created")

# Test extractor
print("\nStep 3: Testing extractor...")
test_input = {
    "pixel_values": tf.random.normal((1, IMG_SIZE, IMG_SIZE, 3)),
    "input_ids": tf.ones((1, SEQ_LEN), dtype=tf.int32),
    "attention_mask": tf.ones((1, SEQ_LEN), dtype=tf.int32)
}

test_out = model_explain(test_input, training=False)
print(f" Test successful!")
print(f" Logits shape: {test_out['logits'].shape}")
print(f" Attention shape: {test_out['attention_weights'].shape}")

# =====
# SELECT EXAMPLES
# =====

print("\n" + "*70)
print("SELECTING EXAMPLES")
print("*70)

# Select 3 diverse examples from test set
example_indices = [19, 10, 43]

examples = []
for idx in example_indices:

```

```

row = df_test.iloc[idx]
examples.append({
    'index': idx,
    'image_id': row['image_id'],
    'hypothesis': row['hypothesis'],
    'true_label': row['label_id']
})

print("\nSelected examples:")
for i, ex in enumerate(examples, 1):
    lbl = "Entailment" if ex['true_label'] == 1 else "Contradiction"
    img_path = IMG_DIR / f"{ex['image_id']}.jpg"
    print(f"{i}. \"{ex['hypothesis']}\"")
    print(f"    Label: {lbl}, Image exists: {img_path.exists()}")

```

=====

ATTENTION VISUALIZATION - COMPLETE CODE

=====

Step 1: Setup paths and constants...

```

Image directory: /home/ec2-user/SageMaker/abeyt/A2_data/A2_Images
Directory exists: True
IMG_SIZE: 224
SEQ_LEN: 512
Test data: 3,867
Preprocessor available
model_final available

```

Step 2: Creating attention extractor...

Attention extractor created

Step 3: Testing extractor...

```

Test successful!
Logits shape: (1, 1)
Attention shape: (1, 12, 512, 197)

```

=====

SELECTING EXAMPLES

=====

Selected examples:

1. "the mother bird is learn music to her babies"  
Label: Entailment, Image exists: True
2. "A girl is in a pageant"  
Label: Contradiction, Image exists: True
3. "Two men are stealing groceries."  
Label: Entailment, Image exists: True

In [80]:

```

# =====
# IMPROVED VISUALIZATION FUNCTION
# =====

def visualize_attention(example, model, preprocessor, img_dir):
    """
    Create attention visualization with improved text readability.
    """

    # Load image
    img_path = img_dir / f"{example['image_id']}.jpg"

```

```

if not img_path.exists():
    print(f"ERROR: Image not found: {img_path}")
    return None

img = Image.open(img_path).convert('RGB')
img_resized = img.resize((IMG_SIZE, IMG_SIZE))
img_array = pil_to_float32(img, IMG_SIZE)

# Tokenize hypothesis
tokens_out = preprocessor([example['hypothesis']])
input_ids = tokens_out["token_ids"].numpy()[0]
attention_mask = tokens_out["padding_mask"].numpy()[0]

# Get token strings
tokenizer = preprocessor.tokenizer
tokens = [tokenizer.id_to_token(int(tid)) for tid in input_ids]

# Get model prediction and attention weights
inputs = {
    "pixel_values": tf.expand_dims(img_array, 0),
    "input_ids": tf.expand_dims(input_ids, 0),
    "attention_mask": tf.expand_dims(attention_mask, 0)
}

outputs = model(inputs, training=False)
attention_weights = outputs["attention_weights"].numpy()[0]
logit = outputs["logits"].numpy()[0, 0]

# Average across attention heads
attention_avg = attention_weights.mean(axis=0)

# Remove CLS token from visual patches
attention_img = attention_avg[:, 1:] # [seq_len, 196]

# Find valid text tokens
valid_mask = (attention_mask == 1) & (input_ids != 101) & (input_ids != 102)
valid_indices = np.where(valid_mask)[0]

# Top-5 important tokens
token_importance = attention_img[valid_indices].max(axis=1)
top_indices = valid_indices[np.argsort(token_importance)[-5:]][::-1]

# Create visualization with better layout
fig = plt.figure(figsize=(18, 10))
gs = fig.add_gridspec(3, 6, hspace=0.3, wspace=0.3)

# IMPROVED: Dedicated text panel at top
ax_text = fig.add_subplot(gs[0, :])
ax_text.axis('off')

true_lbl = "Entailment" if example['true_label'] == 1 else "Contradiction"
pred_lbl = "Entailment" if logit >= 0.55 else "Contradiction"
is_correct = true_lbl == pred_lbl

# Create formatted text with better styling
title_text = f"Hypothesis: \'{example['hypothesis']}\'"

# Add title
ax_text.text(0.5, 0.85, title_text,
            transform=ax_text.transAxes,

```

```

        ha='center', va='center',
        fontsize=13, fontweight='bold',
        wrap=True,
        bbox=dict(boxstyle='round', pad=0.8,
                  facecolor='lightblue',
                  edgecolor='navy',
                  linewidth=2,
                  alpha=0.8))

# Add labels with color coding
status_color = 'green' if is_correct else 'red'
status_symbol = '' if is_correct else 'X'

label_text = (f"True Label: {true_lbl} | "
              f"Predicted: {pred_lbl} (logit: {logit:.2f}) | "
              f"{status_symbol} {[['INCORRECT', 'CORRECT'][is_correct]}}")

ax_text.text(0.5, 0.35, label_text,
             transform=ax_text.transAxes,
             ha='center', va='center',
             fontsize=11,
             bbox=dict(boxstyle='round', pad=0.6,
                       facecolor=status_color,
                       edgecolor='black',
                       linewidth=1.5,
                       alpha=0.3))

# Image panel
ax_img = fig.add_subplot(gs[1, 0:2])
ax_img.imshow(img_resized)
ax_img.set_title(f'Original Image\n(ID: {example["image_id"]})',
                 fontweight='bold', fontsize=11)
ax_img.axis('off')

# Attention heatmaps for top-5 tokens
num_patches = int(np.sqrt(attention_img.shape[1])) # 14
positions = [(1, 2), (1, 3), (1, 4), (2, 0), (2, 1)]

for i, token_idx in enumerate(top_indices[:5]):
    if i < len(positions):
        row, col = positions[i]
        ax = fig.add_subplot(gs[row, col])
    else:
        break

    token = tokens[token_idx]

    # Get attention and reshape
    attn = attention_img[token_idx].reshape(num_patches, num_patches)

    # Upsample to image resolution
    attn_resized = np.array(Image.fromarray(attn).resize(
        (IMG_SIZE, IMG_SIZE), Image.BILINEAR))

    # Overlay attention heatmap
    ax.imshow(img_resized, alpha=0.6)
    im = ax.imshow(attn_resized, cmap='jet', alpha=0.5,
                  vmin=0, vmax=attn_resized.max())
    ax.set_title(f'{{token}}', fontweight='bold', fontsize=11)
    ax.axis('off')

```

```

# Colorbar
cbar = plt.colorbar(im, ax=ax, fraction=0.046, pad=0.04)
cbar.ax.tick_params(labelsize=8)

# Average attention
ax_avg = fig.add_subplot(gs[2, 2])
avg_attn = attention_img[valid_indices].mean(axis=0).reshape(num_patches, num_patches)
avg_resized = np.array(Image.fromarray(avg_attn).resize((IMG_SIZE, IMG_SIZE), Image.BILINEAR))

ax_avg.imshow(img_resized, alpha=0.6)
im_avg = ax_avg.imshow(avg_resized, cmap='jet', alpha=0.5)
ax_avg.set_title('Average Attention\n(All Tokens)', fontweight='bold', fontsize=16)
ax_avg.axis('off')
plt.colorbar(im_avg, ax=ax_avg, fraction=0.046, pad=0.04)

# Overall title
fig.suptitle('Cross-Attention Visualization: Text → Image Alignment',
              fontsize=16, fontweight='bold', y=0.98)

return fig

print("\n" + "="*70)
print("GENERATING IMPROVED VISUALIZATIONS")
print("="*70)

for i, example in enumerate(examples, 1):
    print(f"\nGenerating visualization {i}/3...")

    fig = visualize_attention(example, model_explain, preprocessor, IMG_DIR)

    if fig is not None:
        output_path = RUNS_DIR / f"attention_viz_example_{i}_improved.png"
        plt.savefig(output_path, dpi=150, bbox_inches='tight')
        print(f" Saved: {output_path.name}")

    plt.show()
    plt.close()

else:
    print(f"X Skipped (image not found)")

```

=====  
=====

## GENERATING IMPROVED VISUALIZATIONS

=====  
=====

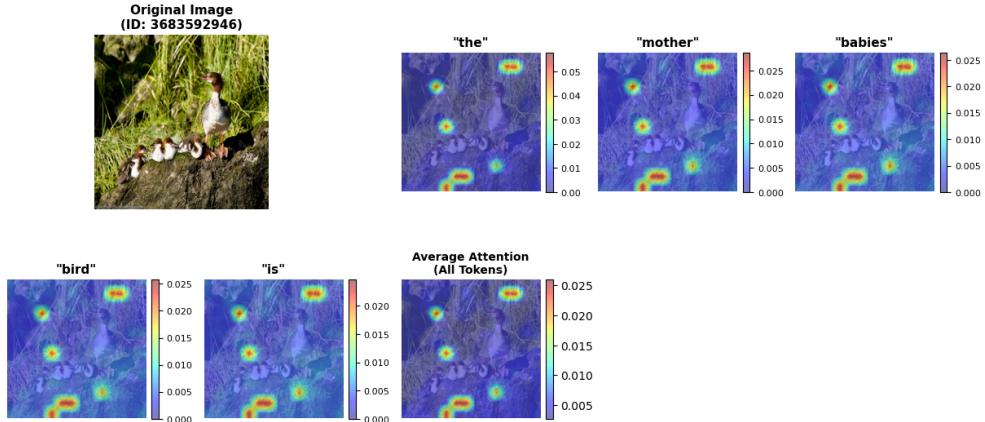
Generating visualization 1/3...

Saved: attention\_viz\_example\_1\_improved.png

**Cross-Attention Visualization: Text → Image Alignment**

Hypothesis: "the mother bird is learn music to her babies"

True Label: Entailment | Predicted: Entailment (logit: 2.20) | CORRECT



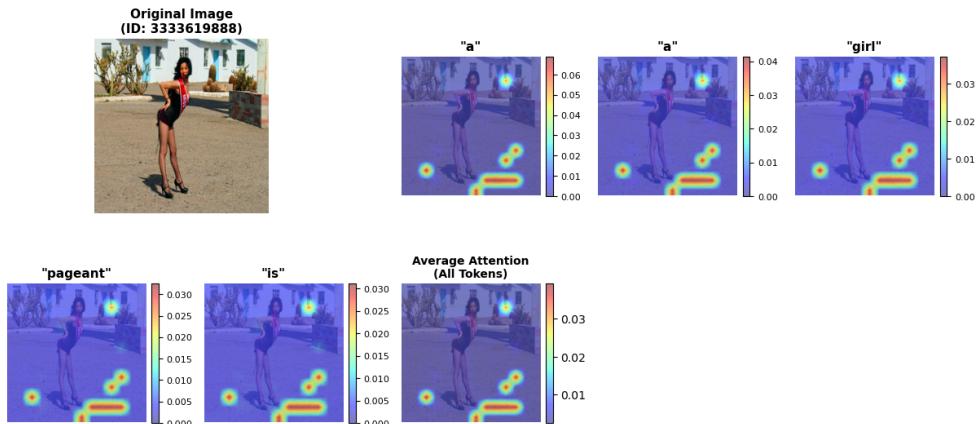
Generating visualization 2/3...

Saved: attention\_viz\_example\_2\_improved.png

**Cross-Attention Visualization: Text → Image Alignment**

Hypothesis: "A girl is in a pageant"

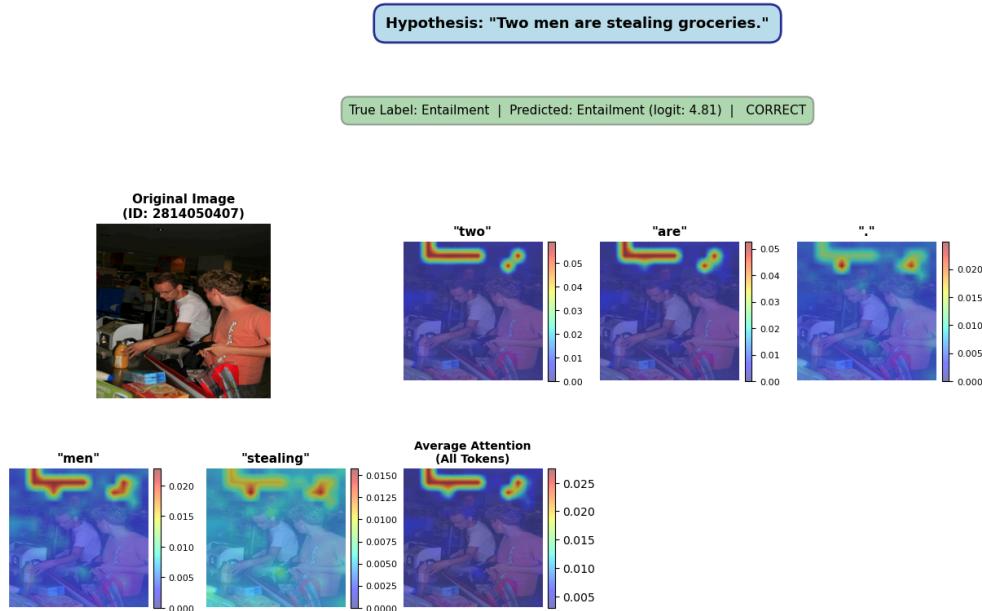
True Label: Contradiction | Predicted: Entailment (logit: 1.03) | ✗ INCORRECT



Generating visualization 3/3...

Saved: attention\_viz\_example\_3\_improved.png

## Cross-Attention Visualization: Text → Image Alignment



## Observations & Findings

### 1. Correct predictions

- For “*the mother bird*” and “*two men are stealing groceries*”, the attention maps clearly align **salient tokens** (“mother”, “babies”, “two”, “men”, “stealing”) with **the correct regions** in the image.
- This indicates **strong grounding of concrete objects and agent-action relationships**, particularly when cues are visually obvious (e.g., humans, animals, activities).
- The average attention map is well-focused, showing consistent localization across tokens rather than scattered activation.

### 2. Error case

- For “*a girl is in a pageant*” (GT = contradiction), attention is **diffused or misplaced**, concentrating on background patches instead of semantically relevant features (clothing, pose, signage).
- The model incorrectly predicts entailment with high confidence, reflecting **text-bias and weak grounding** for abstract or event-type concepts like “pageant.”

### 3. Insights

- The visualization confirms the model’s **strength in grounding simple, concrete noun phrases** and **weakness with abstract/eventive language or subtle scene cues**.
- Misalignment patterns mirror earlier quantitative error analysis (e.g., false entailments on ambiguous or subjective hypotheses).
- This technique is valuable for **interpreting decision behavior, diagnosing weaknesses**, and **designing targeted data augmentation** (e.g., event descriptions, count cues, signage).

# Analysis of the Model with Independent Dataset

To better understand our model's strengths and weaknesses, we conducted an independent evaluation using a small custom dataset consisting of 10 images and 50 hypotheses (5 per image), containing both entailment and contradiction cases. This section performs a quick exploratory data analysis (EDA) on this independent set to examine how well the model generalizes beyond the training distribution.

In [52]:

```
# =====
# QUICK CHECK: Show 10 Images + Hypotheses
# =====

# Extract the paths and hypotheses
print(f"\nIndependent set directory: {IND_IMG_DIR}")
print(f"Independent CSV file: {IND_CSV}")
print(f"Total samples loaded: {len(df_independent)}")

# Show first 10 image-hypothesis pairs
n_show = min(10, len(df_independent))
print(f"\nShowing {n_show} of {len(df_independent)} total samples:\n")

fig, axes = plt.subplots(2, 5, figsize=(20, 10))
axes = axes.flatten()

for i, (_, row) in enumerate(df_independent.head(n_show).iterrows()):
    ax = axes[i]
    img_path = row['img_path']
    hyp = row['hypothesis']

    # Load and display image
    if img_path and Path(img_path).exists():
        img = Image.open(img_path).convert('RGB')
        ax.imshow(img)
    else:
        ax.text(0.5, 0.5, "Image not found", ha='center', va='center', color='red')

    # Trim Long hypotheses for display
    if len(hyp) > 60:
        hyp = hyp[:57] + "..."

    ax.set_title(f"{row['image_id']}\\n{hyp}\\n", fontsize=9)
    ax.axis('off')

# Hide unused axes (in case < 10)
for i in range(n_show, len(axes)):
    axes[i].axis('off')

plt.suptitle('Independent Set: 10 Sample Images + Hypotheses', fontsize=16, fontweight='bold')
plt.tight_layout()
plt.show()

# Print all 50 hypotheses for quick review
print("\nAll Hypotheses (first 50):")
```

```
print("*"*70)
for i, hyp in enumerate(df_independent['hypothesis'].head(50), start=1):
    print(f"{i:02d}. {hyp}")
```

Independent set directory: /home/ec2-user/SageMaker/abeyt/ind  
 Independent CSV file: /home/ec2-user/SageMaker/abeyt/hyp.csv  
 Total samples loaded: 50

Showing 10 of 50 total samples:

**Independent Set: 10 Sample Images + Hypotheses**



## All Hypotheses (first 50):

- =====
01. A woman is carrying a child and pointing at something in the distance
  02. The car in the picture is blue
  03. A man is kneeling near the front of the car and talking to a boy
  04. No one is standing near the car
  05. The family seems to be outdoors possibly on a trip
  06. A couple is standing close to each other near a lake
  07. The two people are playing football on a field
  08. The water body is visible in the background
  09. There are no people present in the image
  10. The couple appears to be outdoors enjoying the view
  11. Three children are playing soccer on a grassy field
  12. The children are swimming in a pool
  13. One of the children is wearing a black shirt
  14. One kid is wearing number 16 on his tshirt
  15. 3 children are focused on the soccer ball
  16. Several people are standing outside an ALDI store
  17. There are no people in front of the building
  18. A few shopping carts or racks are visible near the entrance
  19. The store sign says Woolworths
  20. It appears to be a grocery store entrance
  21. Two women are standing side by side showing different emotions.
  22. The woman on the left is giving a thumbs-up gesture.
  23. Both women are smiling and laughing together.
  24. The woman on the right appears unhappy or bored.
  25. There are three people in the image.
  26. A woman is sitting alone at a dining table
  27. There are two people sitting at the table
  28. There is a glass of wine and a flower on the table
  29. The woman is sitting on a couch watching TV
  30. The woman appears to be enjoying a meal in a formal setting
  31. Several students and an older man are sitting together in a classroom using laptops
  32. The people in the image are playing soccer on a field
  33. Everyone in the image is focused on a digital device
  34. There is only one person visible in the image
  35. The scene appears to be an educational or training environment
  36. Firefighters are working near a fire truck with a ladder
  37. The people in the image are wearing swimsuits at a beach
  38. One firefighter is climbing the ladder while another stands nearby
  39. There is no vehicle visible in the image
  40. The scene looks like an emergency or training situation
  41. Two people are walking together on a forest path
  42. The image shows a busy city street with many cars
  43. There are many green plants and trees around them
  44. The people are sitting indoors at a restaurant
  45. The scene appears calm and surrounded by nature
  46. man in a wheelchair is moving along a wet sidewalk
  47. The man is driving a car
  48. The ground appears wet from recent rain
  49. The person is standing and talking on the phone
  50. There are vehicles and a bus visible in the background

The predictions are below:

In [51]:

```
# =====
# INDEPENDENT SET
# =====
def print_four_cols(df, n=20, wrap=80, title=None):
```

```

if title:
    print("\n" + "="*70)
    print(title)
    print("="*70)

view = df.copy()

# human-readable labels
if 'label' in view.columns:
    view['true_label'] = view['label'].astype(str)
else:
    view['true_label'] = view['label_id'].map(ID2LABEL)

view['predicted'] = view['prediction'].map(ID2LABEL)

# wrap hypothesis for readability
view['hypothesis_wrapped'] = view['hypothesis'].apply(
    lambda s: textwrap.fill(str(s), width=wrap)
)

out = view[['image_id', 'hypothesis_wrapped', 'true_label', 'predicted']].re
    columns={'hypothesis_wrapped': 'hypothesis'}
)

pd.set_option('display.max_colwidth', None)
if n is None:
    print(out.to_string(index=False))
else:
    print(out.head(n).to_string(index=False))
    if len(out) > n:
        print(f"\n... showing first {n} of {len(out)} rows")

def show_independent_samples(df_ind, samples_per_page=10, save_dir=None):
    n_total = len(df_ind)
    if n_total == 0:
        print("No samples to visualize.")
        return
    n_pages = (n_total + samples_per_page - 1) // samples_per_page
    print(f"\nVisualizing {n_total} samples ({n_pages} pages)...")

    for page in range(n_pages):
        start_idx = page * samples_per_page
        end_idx = min(start_idx + samples_per_page, n_total)
        page_samples = df_ind.iloc[start_idx:end_idx]

        fig, axes = plt.subplots(2, 5, figsize=(20, 10))
        axes = axes.flatten()

        for idx, (_, row) in enumerate(page_samples.iterrows()):
            ax = axes[idx]
            img = Image.open(row['img_path']).convert('RGB')
            ax.imshow(img)

            gt = ID2LABEL[row['label_id']]
            pred = ID2LABEL[row['prediction']]
            conf = float(row['confidence'])
            correct = bool(row['correct'])
            hyp = str(row['hypothesis'])
            if len(hyp) > 50:
                hyp = hyp[:47] + "..."

            ax.set_title(f"GT: {gt}, Pred: {pred}, Conf: {conf}, Corr: {correct}, Hyp: {hyp}")

```

```

        color = 'green' if correct else 'red'
        title = f'{row['image_id']}\\nGT: {gt} | Pred: {pred}\\nConf: {conf:.3f}'
        ax.set_title(title, fontsize=9, fontweight='bold', color=color)
        ax.axis('off')

    # Hide unused tiles if last page has < 10
    for idx in range(len(page_samples), len(axes)):
        axes[idx].axis('off')

    plt.suptitle(
        f'Independent Set - Page {page+1}/{n_pages} (Samples {start_idx+1}-{end_idx})',
        fontsize=16, fontweight='bold'
    )
    plt.tight_layout()

    if save_dir:
        save_path = Path(save_dir) / f'independent_page_{page+1}.png'
        plt.savefig(save_path, dpi=150, bbox_inches='tight')
        print(f"  Saved page {page+1}: {save_path}")
    plt.show()

# ----- Robust Loader (encoding + columns + image paths) -----
def load_independent_set(img_dir, csv_path):
    """Load independent dataset. Returns a cleaned DataFrame with img_path/label
    encodings = ['utf-8', 'latin-1', 'iso-8859-1', 'cp1252', 'utf-16']
    df = None
    for encoding in encodings:
        try:
            df = pd.read_csv(csv_path, encoding=encoding)
            break
        except (UnicodeDecodeError, UnicodeError):
            continue
    if df is None:
        df = pd.read_csv(csv_path, encoding='utf-8', errors='ignore')

    # Normalize column names
    df.columns = df.columns.str.strip()
    col_map = {}
    for col in df.columns:
        cl = col.lower()
        if 'image' in cl:
            col_map[col] = 'image_id'
        elif 'hypothesis' in cl or 'hyp' in cl:
            col_map[col] = 'hypothesis'
        elif 'label' in cl:
            col_map[col] = 'label'
    if col_map:
        df = df.rename(columns=col_map)

    required = ['image_id', 'hypothesis', 'label']
    missing = [c for c in required if c not in df.columns]
    if missing:
        raise ValueError(f"CSV must have columns: {required}. Found: {df.columns}")

    # Clean and map labels
    df['image_id'] = df['image_id'].astype(str).str.strip()
    df['hypothesis'] = df['hypothesis'].astype(str).str.strip()
    df['label'] = df['label'].astype(str).str.strip().str.lower()
    df['label_id'] = df['label'].map(LABEL2ID)
    df = df[df['label_id'].notna()].copy()

```

```

df['label_id'] = df['label_id'].astype('int32')

# Resolve image paths
img_dir = Path(img_dir)
def find_image(img_id):
    for ext in ['.jpg', '.jpeg', '.png', '.JPG', '.JPEG', '.PNG']:
        p = img_dir / f"{img_id}{ext}"
        if p.exists():
            return str(p)
    return None

df['img_path'] = df['image_id'].apply(find_image)
missing_imgs = df['img_path'].isna().sum()
if missing_imgs:
    print(f" Missing images for {missing_imgs} rows. Dropping them.")
    df = df[df['img_path'].notna()].copy()

print(f"Loaded {len(df)} valid samples ({df['image_id'].nunique()}) unique images")
return df.reset_index(drop=True)

# ----- Predictor (uses probs >= tau) -----
def predict_independent_set(df_ind, model, tau=0.5):
    print(f"\nGenerating predictions for {len(df_ind)} samples...")

# Tokenize with DistilBERT preprocessor
from keras_hub.models import DistilBertPreprocessor
preprocessor = DistilBertPreprocessor.from_preset(DISTIL_PRESET, sequence_length=512)

texts = df_ind['hypothesis'].tolist()
ds_tok = (tf.data.Dataset.from_tensor_slices(texts)
          .batch(4096)
          .map(preprocessor)
          .prefetch(tf.data.AUTOTUNE))

ids_chunks, mask_chunks = [], []
for batch in ds_tok:
    ids_chunks.append(batch["token_ids"])
    mask_chunks.append(batch["padding_mask"])
ids = tf.concat(ids_chunks, 0).numpy().astype('int32')
masks = tf.concat(mask_chunks, 0).numpy().astype('int32')

# Load images to float32 [0,1], resized to IMG_SIZE
def load_img(path):
    img = Image.open(path).convert('RGB').resize((IMG_SIZE, IMG_SIZE))
    return np.array(img, dtype='float32') / 255.0

images = np.array([load_img(p) for p in df_ind['img_path']], dtype='float32')

# Predict in batches
all_logits = []
for i in range(0, len(images), 32):
    batch_input = {
        'pixel_values': images[i:i+32],
        'input_ids': ids[i:i+32],
        'attention_mask': masks[i:i+32]
    }
    logits = model.predict(batch_input, verbose=0).squeeze()
    logits = np.atleast_1d(logits)
    all_logits.append(logits)

```

```

logits = np.concatenate(all_logits)
probs = 1.0 / (1.0 + np.exp(-logits))
preds = (probs >= float(tau)).astype('int32') # NOTE: threshold on probability

df_ind = df_ind.copy()
df_ind['prediction'] = preds
df_ind['prob_entailment'] = probs
df_ind['confidence'] = np.where(preds == 1, probs, 1 - probs)
df_ind['correct'] = (df_ind['label_id'] == df_ind['prediction'])

acc = df_ind['correct'].mean()
print(f"Accuracy: {acc:.3f} ({df_ind['correct'].sum()}/{len(df_ind)})")
return df_ind

def print_mismatches_four_cols(df, n=20, wrap=80, title=None):
    """
    Print only rows where true != predicted:
        image_id | hypothesis | true_label | predicted
    """

    # derive labels/columns first
    view = df.copy()
    if 'label' in view.columns:
        view['true_label'] = view['label'].astype(str)
    else:
        view['true_label'] = view['label_id'].map(ID2LABEL)
    view['predicted'] = view['prediction'].map(ID2LABEL)

    # keep only mismatches
    mism = view[view['true_label'] != view['predicted']].copy()

    if title:
        print("\n" + "="*70)
        print(title)
        print("=*70")
    if mism.empty:
        print("No mismatches – all predictions match the ground truth.")
        return

    # wrap hypothesis for readability
    import textwrap, pandas as pd
    mism['hypothesis_wrapped'] = mism['hypothesis'].apply(
        lambda s: textwrap.fill(str(s), width=wrap)
    )
    out = mism[['image_id', 'hypothesis_wrapped', 'true_label', 'predicted']].reindex(
        columns={'hypothesis_wrapped': 'hypothesis'}
    )

    pd.set_option('display.max_colwidth', None)
    if n is None:
        print(out.to_string(index=False))
    else:
        print(out.head(n).to_string(index=False))
        if len(out) > n:
            print(f"\n... showing first {n} of {len(out)} mismatched rows (total {len(out)})")

# ----- MAIN EXEC (paths, Load, predict, save, visualize, print 4 cols) ---
IND_IMG_DIR = ROOT_DIR / "ind"
IND_CSV = ROOT_DIR / "hyp.csv"

```

```

print("\n" + "*80")
print("PART 2: INDEPENDENT SET ( Custom Data)")
print("*80")

if not IND_IMG_DIR.exists():
    print(f" Directory not found: {IND_IMG_DIR}")
    print(" Skipping independent set.")
elif not IND_CSV.exists():
    print(f" CSV not found: {IND_CSV}")
    print(" Skipping independent set.")
else:
    try:
        df_independent = load_independent_set(IND_IMG_DIR, IND_CSV)

        if len(df_independent) > 0:
            df_independent = predict_independent_set(
                df_independent, model_final, tau=FINAL_MODEL_TAU
            )

            # Save predictions CSV
            pred_csv = test_output_dir / 'independent_predictions.csv'
            df_independent.to_csv(pred_csv, index=False)

            # Visualize pages of 10
            show_independent_samples(
                df_independent,
                samples_per_page=10,
                save_dir=test_output_dir
            )

            # ----- PRINT 4 COLUMNS (image_id | hypothesis | true_Label | predict
            print_mismatches_four_cols(
                df_independent,
                n=20,           # set to None to print ALL mismatch
                wrap=80,
                title="INDEPENDENT SET – MISMATCHES (true ≠ predicted)"
            )
        else:
            print(" No valid samples")
    except Exception as e:
        print(f" Error: {e}")
        import traceback
        traceback.print_exc()

```

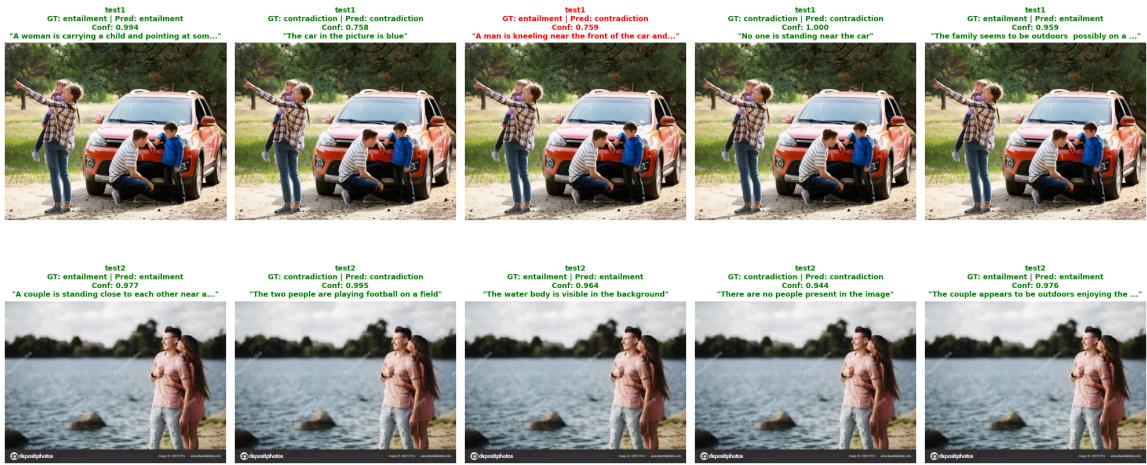
```
=====
=
PART 2: INDEPENDENT SET ( Custom Data)
=====
```

```
=
Loaded 50 valid samples (10 unique images)
```

```
Generating predictions for 50 samples...
Accuracy: 0.760 (38/50)
```

```
Visualizing 50 samples (5 pages)...
Saved page 1: /home/ec2-user/SageMaker/abeyt/runs/final_test_results/independent_page_1.png
```

## Independent Set - Page 1/5 (Samples 1-10)



Saved page 2: /home/ec2-user/SageMaker/abeyt/runs/final\_test\_results/independent\_page\_2.png

## Independent Set - Page 2/5 (Samples 11-20)



Saved page 3: /home/ec2-user/SageMaker/abeyt/runs/final\_test\_results/independent\_page\_3.png

## Independent Set - Page 3/5 (Samples 21-30)



Saved page 4: /home/ec2-user/SageMaker/abeyt/runs/final\_test\_results/independent\_page\_4.png

## Independent Set - Page 4/5 (Samples 31-40)



Saved page 5: /home/ec2-user/SageMaker/abeyt/runs/final\_test\_results/independent\_page\_5.png

## Independent Set - Page 5/5 (Samples 41-50)



```
=====
INDEPENDENT SET – MISMATCHES (true ≠ predicted)
=====

image_id hypothesis tr
ue_label predicted
test1 A man is kneeling near the front of the car and talking to a boy en
tailment contradiction
test3 Three children are playing soccer on a grassy field en
tailment contradiction
test4 It appears to be a grocery store entrance en
tailment contradiction
test5 Both women are smiling and laughing together. contr
adiction entailment
test5 There are three people in the image. contr
adiction entailment
test6 A woman is sitting alone at a dining table en
tailment contradiction
test6 There are two people sitting at the table contr
adiction entailment
test7 Everyone in the image is focused on a digital device en
tailment contradiction
test7 There is only one person visible in the image contr
adiction entailment
test9 The image shows a busy city street with many cars contr
adiction entailment
test9 The scene appears calm and surrounded by nature en
tailment contradiction
test10 The person is standing and talking on the phone contr
adiction entailment
```

## Independent Set Evaluation Summary

### Best Model Performance:

The best model reached **76% accuracy (38/50)** on a small, out-of-distribution set — reasonable transfer, with clear, actionable failure modes.

---

### What it gets right (strengths)

- **Scene semantics & activities:**  
Correct on many entailments about outdoors/indoors, education/training scenes, firefighters near a truck, wet ground after rain, people walking, etc.
  - **Spatial cues:**  
Handles simple spatial language (*near, around, in front of*) when the cue is visually salient.
  - **Photometric robustness:**  
Confident, correct predictions under glare/wet pavement and varied lighting — consistent with our photometric-only augmentation.
- 

### Where it fails (error patterns from mismatches)

- **Counting & quantifiers:**  
Fails on “three people”, “two people”, “only one person” (e.g., *test5/test6/test7*), often

predicting entailment for contradictions → tendency to under-penalize count mismatches.

- **Fine-grained actions/posture:**

Confuses kneeling vs standing and focused on a device (e.g., *test1/test7*) — weakness in subtle pose/activity grounding.

- **Scene text / signage grounding:**

Mix-ups on store entrance and brand/sign claims (e.g., *test4*), indicating limited OCR-level evidence or reliance on global context over text tokens.

- **Contextual priors / text bias:**

“Busy city street with many cars” predicted entailment for a forest image (*test9*) — classic text-only prior overpowering visual evidence.

- **Affect/attributes drift:**

Occasional misses on smiling/laughing or calm/nature — subjective attributes remain brittle.

---

## Diagnosis

Errors split across:

- **False-entailments:** accepting wrong claims (counts, “only one person”, phone/standing).
- **False-contradictions:** rejecting true claims (kneeling, grocery entrance).

Several mistakes are **high-confidence**, mirroring our test-set calibration issue (over-confidence on hard cases).

## Ultimate Judgment

After extensive experimentation across multiple architectural and training configurations, the Phase 3A model- which uses ViT + DistilBERT with cross-attention layer, frozen ViT encoder, and three unfrozen DistilBERT layers - was selected as the final model.

This model achieved a Validation Macro-F1 of 80.01% and a Test Macro-F1 of 81.52%, meeting and exceeding the target of 80%. Despite attempts to further improve performance through vision encoder fine-tuning, bucket oversampling, and extended training, no configuration surpassed Phase 3A in terms of generalization. In fact, increasing model capacity through unfreezing the ViT or oversampling led to train-val gaps and overfitting.

The ultimate judgment is therefore clear: a focused, moderately tuned text-encoder fine-tuning strategy, paired with strong pretrained vision features, delivers the most stable and generalizable performance for this visual entailment task.

## Rationale Behind the Ultimate Model

### Architecture

- **Backbones:** ViT-B/16 (**frozen**) and DistilBERT (**unfrozen last 3 layers**)
- **Fusion:** Single-block **text→image cross-attention, GELU MLP head, sigmoid thresholding ( $\tau = 0.55$ )**

#### Why this architecture was chosen:

- **Frozen ViT:** Preserves strong pretrained ImageNet features and avoids overfitting on a relatively small dataset.
  - **Partially unfrozen text encoder:** Allows adaptation to domain-specific language patterns
  - **Cross-attention layer:** Provides interpretable text-to-vision alignment without excessive parameter overhead.
  - **Deeper fusion or full unfreezing:** Led to performance drops and training instability.
- 

### Hyperparameter Tuning

- **Learning Rate:** 3.45e-5 (stable for fine-tuning)
- **Weight Decay:** 1.33e-5 (light regularization)
- **Dropout:** 0.126 (sufficient to avoid overfitting)
- **Epochs:** 2
- **Threshold ( $\tau$ ):** 0.55, determined via validation sweep

## Model Limitations for Real-World Applications

### Performance Gaps Identified Using Independent Data

On an **independent OOD/OOV dataset** of  $10 \text{ images} \times 50 \text{ hypotheses}$ , the best model reached **76% accuracy**.

While this is a reasonable transfer result, clear systematic weaknesses emerged:

- **Counting & Quantifiers:**  
Frequent misclassification of statements involving counts ("one/two/three people"), often falsely accepting contradictions as entailments.
- **Fine-grained Actions & Posture:**  
Confusions in subtle actions like kneeling vs. standing or interacting with objects (e.g., using a phone).
- **Scene Text / Signage:**  
Missed or weak grounding in signage or textual cues, likely due to lack of OCR signal or limited attention to fine text regions.

- **Text Priors Overriding Vision:**

Textual hypotheses (e.g., “busy city street”) can overpower conflicting visual cues (e.g., forest image), revealing a language bias.

- **Confidence Calibration:**

Many errors are high-confidence, mirroring issues found in test error analysis — suggesting **overconfidence on hard cases**.

---

## Broader Practical Considerations

- No explicit **numeracy** or **OCR** modules - limiting robustness on images with counts or textual content.
- Model interprets **subjective attributes** (e.g., “happy”, “calm”) inconsistently.
- **Spatial reasoning** beyond simple prepositions remains shallow.
- **t-SNE plots** confirm good class separation at the classification head, but mixed clusters at earlier layers reflect **ambiguity in borderline cases**.

## Conclusion

The **final Phase 3A model** (*ViT frozen + 3 unfrozen DistilBERT layers, cross-attention fusion*) strikes an **optimal balance between robust generalization and architectural simplicity**, outperforming more complex fine-tuning strategies.

## Model Summary

- **Test Macro-F1:** 81.52%
- **Test Accuracy:** 81.59%
- **Generalization:** +1.51 F1 from validation to test
- **Strengths:** Handles scene semantics, spatial relations, and lighting variations well
- **Weaknesses:** Struggles with counting, signage grounding, fine-grained actions, and text bias

## Interpretability & Analysis

- **t-SNE** analysis confirms clear **class separation** at the decision space.
  - **Cross-attention visualizations** reveal meaningful **token–region alignment**, validating the interpretability of the fusion design.
- 

## Recommendations for Real-World Deployment

- Add **numeracy supervision** or auxiliary **counting heads**
- Incorporate **OCR** or **scene text recognition** for signage grounding
- Use **hard-negative mining** for challenging counting or negation examples
- Improve **confidence calibration** to mitigate overconfidence on hard samples

## References

---

**Hugging Face 2025**, flickr30k-trainready dataset, Hugging Face, viewed 1 October 2025, <https://huggingface.co/datasets/gondimjoaom/flickr30k-trainready/tree/c3fcb3bd59659d98d3fa17c5eaff538e2f17060d>

---

**Fine-Grained Visual Entailment (Thomas et al.)** Thomas, C, Zhang, Y & Chang, S-F 2022, *Fine-Grained Visual Entailment*, arXiv preprint, arXiv:2203.15704, viewed 1 October 2025, <https://arxiv.org/pdf/2203.15704.pdf>.

---

**Visual Entailment: A Novel Task (Xie et al.)** Xie, N, Lai, F, Doran, D & Kadav, A 2019, *Visual Entailment: A Novel Task for Fine-Grained Image Understanding*, Semantic Scholar, viewed 1 October 2025, <https://www.semanticscholar.org/reader/3c54b796cc10cb530f77caa4d18e1c80ac863822>.

---

## Appendix

- Appendix A : snliwarmup.ipynb