**Emanuele Vivoli,**
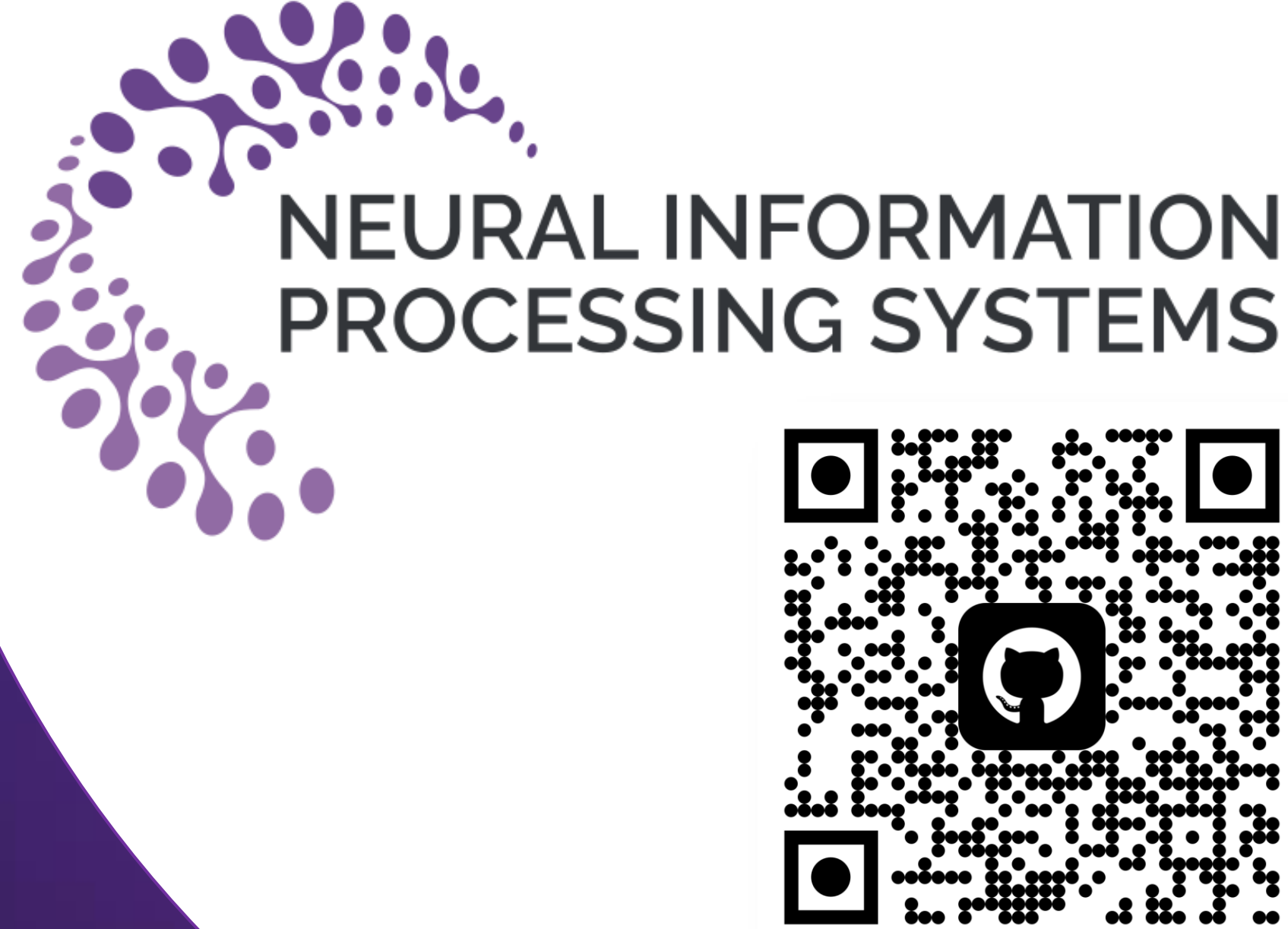Marco Bertini,
Dimosthenis Karatzas

# CoMix : A comprehensive Benchmark for Multi-Task Comic Understanding

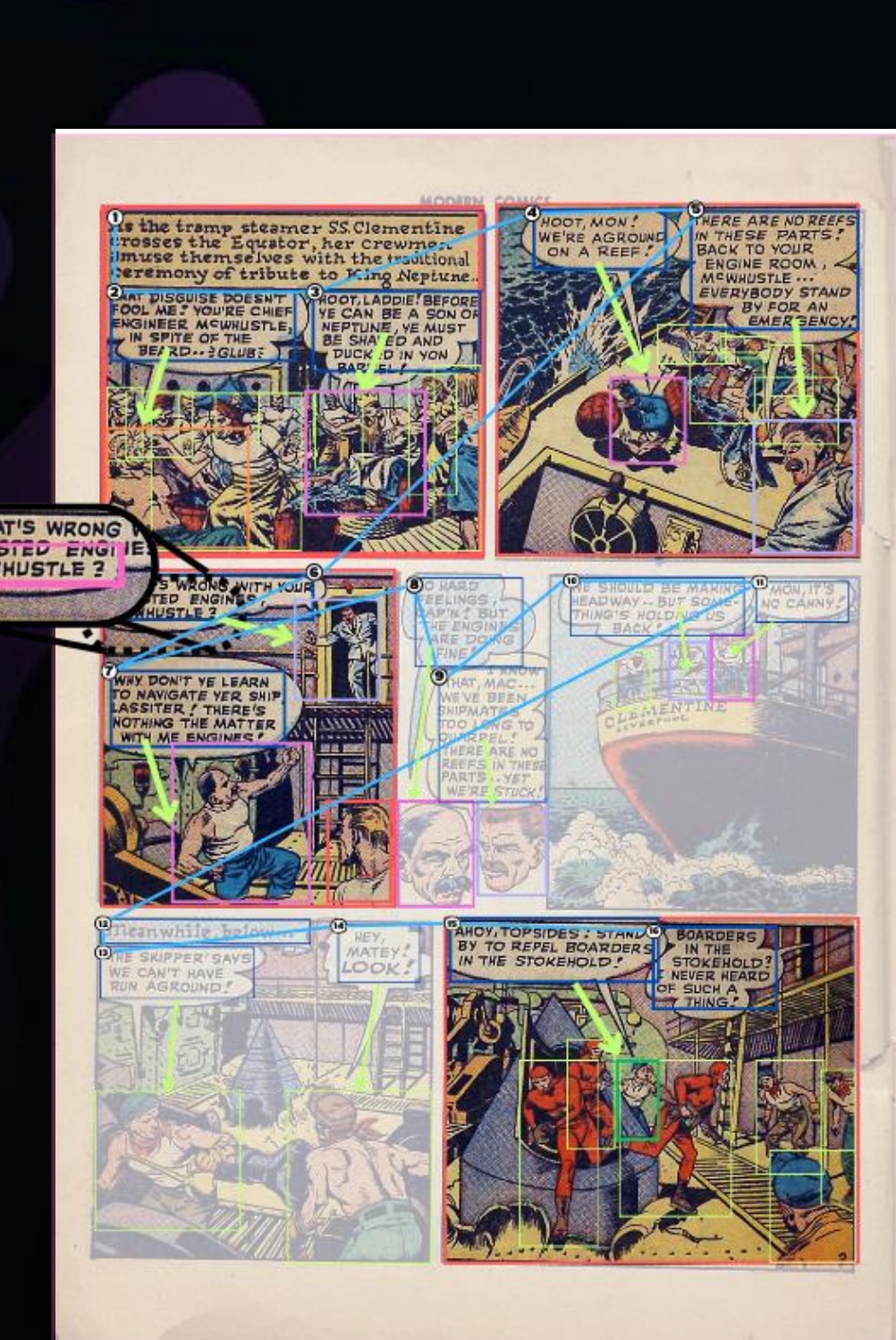NEURAL INFORMATION PROCESSING SYSTEMS

CVC — Computer Vision Center

micc

## Motivation

Comics datasets only tackle **simple tasks**, are **small** or of a **single style** (manga vs. comics)

| Dataset | Release | Avail | Tasks | Years | Style | Books | Pages |
|---|---|---|---|---|---|---|---|
| eBDtheque | 2013 | ✓ | d,t2c | 1905-2012 | mix | 28 | 100 |
| COMICS | 2017 | ✓ | c | 1938-1954 | comics | 3948 | 198k |
| GCN | 2017 | ✗ | d,t2c | 1978-2013 | comics | *253 | *38k |
| DCM772 | 2018 | ✓ | d | 1938-1954 | comics | 27 | 772 |
| Manga109 | 2018 | ✓ | d,t2c,c2c | 1970-2010 | manga | 109 | 10k |
| BCBId | 2022 | ✓ | - | - | bangla | 64 | 3k |
| VLRC | 2023 | ✗ | - | 1940-now | - | *376 | *7k |
| PopManga | 2024 | ✓ | d,t2c,c2c | 2010-2023 | manga | 25 | 1.8k |
| *CoMix* (our) | 2024 | ✓ | d,t2c,c2c,N,D | 1938-2023 | mix | 100 | 3.8k |

We create a densely annotated dataset: *CoMix*

| Legend | | |
|---|---|---|
| d: Detection | t2c: Text-to-Character | N: Character naming |
| c: Classification | c2c: Char-to-Char | D: Dialog generation |

## Design choices

Comics    DCM    eBDtheque    PopManga

Datasets — Before / After

Tasks — Before / After

eBDtheque 2.3%, DCM 21.5%, popmanga 50%, COMICS 26.3%

## Annotations

We improve the quality of existing annotations

BEFORE    AFTER

And provide annotations for new tasks

## Tasks

**Object Detection** — mAP R@k
Panels    Characters    Textboxes    Faces

**Speaker id.** — R@#text

**Character Re-Id** — AMI NMI

**Reading order** — Edit distance

**Character Naming** — ANLS
Narrator, Sailor 1, McWhustle, Captain Matey, Sailor 2, Sailor 3, Sailor 4

**Dialog generation** — Hybrid Dialog Score

**Narrator:** "As the tramp steamer SS. Clementine crosses the Equator [...]"
**Sailor 1:** "THAT DISGUISE DOESN'T FOOL ME! YOU'RE CHIEF ENGINEER [...]"
**McWhustle:** "HOOT, LADDIE! BEFORE YE CAN BE A SON OF NEPTUNE, YE [...]"
**McWhustle:** "HOOT, MON! WE'RE AGROUND ON A REEF!"
**Captain:** "THERE ARE NO REEFS IN THESE PARTS! BACK TO YOUR ENGINE [...]"
**Captain:** "WHAT'S WRONG WITH YOUR BLASTED ENGINES, MCWHUSTLE ?"
**McWhustle:** "WHY DON'T YE LEARN TO NAVIGATE YER SHIP, LASSITER! [...]"
**McWhustle:** "NO HARD FEELINGS, CAP'N! BUT THE ENGINES ARE DOING FINE!"
**Captain:** "I KNOW THAT, MAC... WE'VE BEEN SHIPMATES TOO LONG TO [...]"
**Captain:** "WE SHOULD BE MAKING HEADWAY... BUT SOMETHING'S [...]"
**McWhustle:** "MON, IT'S NO CANNY!"
**Narrator:** "Meanwhile, below..."
**Matey:** "'THE SKIPPER' SAYS WE CAN'T HAVE RUN AGROUND!"
Sailor 2: "HEY, MATEY! LOOK!"
**Sailor 3:** "AHOY, TOPSIDES! STAND BY TO REPEL BOARDERS IN THE [...]"
Sailor 4: "BOARDERS IN THE STOKEHOLD? I NEVER HEARD OF SUCH A THING."

## HDS Metric

Hungarian Matching with Edit distance

GT    pred

**Ground Truth**
**Uncle Sam:** "GRR!! I SNAP YOUR NECK LIKE PIGEON!"
**Char I:** "NO! GASPS NO! DON'T! AGGG-GGI"
**Iron Ace:** "WELL, THERE'S A MAN WHO DOES A GOOD [...]"
**Iron Ace:** "I'M GETTING MY USUAL GOOD-BYE... I DON'T [...]"
**Iron Ace:** "NOW, IF THEY ONLY KEEP DOING THE [...]"
**Radio Operator:** "PATROL PLANES ATTENTION... [...]"
**Iron Ace:** "OH-OH!! THE WHOLE JAP AIR FORCE IS [...]"
**Iron Ace:** "WELL, COME AND GET IT, BOYS... BUT IT'S THE [...]"

**GPT4**
**US Soldier:** "GRR!! I SNAP YOUR NECK LIKE PIGEON!"
**Pilot 1:** "NO! GASPS NO! DON'T! AGGG-GGI"
**Pilot 1:** "WELL, THERE'S A MAN WHO DOES A GOOD [...]"
**Pilot 1:** "NOW, IF THEY ONLY KEEP DOING THE MISSING [...]"
**Pilot 1:** "OH-OH!! THE WHOLE JAP AIR FORCE IS AFTER [...]"
**Pilot 1:** "I'M GETTING MY USUAL GOOD-BYE... I DON'T [...]"
**Radio Operator:** "PATROL PLANES ATTENTION... [...]"
**Pilot 1:** "WELL COME AND GET IT, BOYS... BUT IT'S THE [...]"

names    sentences    matches

```
Algorithm 2 Hybrid Dialog Score
1: procedure EVALUATETRANSCRIPTION(model_output, ground_truth)
2:     matches ← find optimal matches(model_output, ground_truth)
3:     tot_ed, char_name_score ← 0,0,0
4:     for each (mo, gt) in matches do
5:         edit_dist ← calculate edit distance(mo.text, gt.text)
6:         tot_ed ← tot_ed + edit_dist / len(gt.text)
7:         ans_score ← calculate ANLS (mo.name, gt.name)
8:         char_name_score ← char_name_score + ans_score
9:     end for
10:    tot_ed ← 1 − tot_ed
11:    char_name_score ← char_name_score / len(matches)
12:    return tot_ed, char_name_score
13: end procedure
```

Σ  Edit distance / Length

Σ  ANLS / Length matches

## Benchmarks

We evaluate models on a "per-task" setting as no model can perform all the 6 tasks together:

| Task | Output | Metric | Baseline | Score |
|---|---|---|---|---|
| Object detection | box detection | mAP - R@100 | Magi | 78.6 - 67.9 |
| Speaker identification | object indexes | R@#text | heuristic | 0.68 |
| Character Re-Id | cluster ids | AMI - NMI | DINOv2 | 0.29 - 0.51 |
| Character Naming | names | ANLS | GPT-4 | 47.11 |
| Dialog generation | list of tuples | HDS | GPT-4 | 93.14 |

This is Lele
if you like Comics/Manga and Vision-Language, come and talk to me!

**and... read our SURVEY**