# GE ADJUNCT TEST REPORT: OBJECT-ORIENTED DESIGN AND SCORING FOR MUC-4

George Krupka and Lisa Rau Artificial Intelligence Laboratory GE Research and Development Schenectady, NY 12301 USA E-mail: rau@crd.ge.com Phone: (518) 387 - 5059

#### Abstract

This paper reports on the results of the adjunct test performed by GE for the MUC-4 evaluation of text processing systems. In this test, we evaluated the effect of an object-oriented template design and associated matching conditions on the scores. The results indicate that the current MUC-4 "flat" templade design with cross-references closely approximates a true object-oriented design. However the object-oriented design allows for additional performance data to be calculated, facilitating diagnosis.

## INTRODUCTION

In this adjunct test, we investigate the issues and effect of transforming the MUC-4 template design automatically into an object-oriented design, with associated object-level matching conditions affecting the overall score.

An object is simply a collection of slots that all refer to one item originating in the text. This collection of slots is logically connected due to their implicit reference to one particular filler. In addition, objects may be nested so that one object contains another object, or may be recursive, with one object pointing to another and back again. Finally, there may be multiple instances of any given object.

In MUC-4 is is possible to isolate three distinct levels of objects. The first level contains a STORY object. Attached to a STORY object are INCIDENT objects, each of which contains participant objects; TARGET, INSTRUMENT and PERPETRATOR. The Figure 1 illustrates an object-oriented template design.

The MUC template design has been moving from a flat structure to a more object-oriented structure with cross-references tying together multiple slots into what has been termed a "pseudo-object". Slots tied together with cross-references are pseudo-objects and not true objects because the cross-references are not enforced. For example, consider the following scenario depicted in Figure 2. Although the human target ('MARY') is wrong, the system is allowed partial credit for the other slots, even though they are clearly cross-referenced to different, and equal incorrect, fills.

## MOTIVATION

There are a variety of reasons why an object-oriented design is desirable. First, an object-oriented design is conceptually easier to understand. Instead of a flat listing of all slots of a template, slots pertaining to a single fill are grouped together. Cross-references are no longer needed, as indentation indicates the grouping, so the visual design is cleaner.

Asthetics aside, there are additional performance data that can be obtained when groups of slots are connected as objects. Systems can be scored on how well a given object aligns, in a way analogous to template matching alignment scores (template ID score). Moreover, it is possible to construct an object total that does not consider whether any given object appears in a *matching* template or not. For example in our system, we found that data such as a human target was correctly extracted by the program with all its associated fields, but was put in the wrong template. This resulted in missing and spurious points for all the

# Story

slot: incident

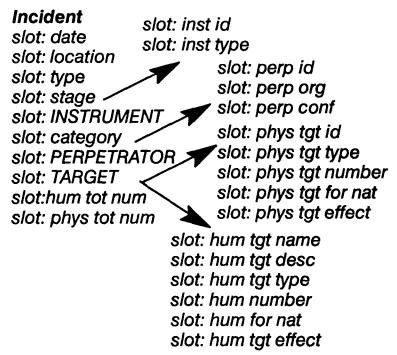


Figure 1: Object-oriented MUC-4 Template Design

fields in the object, although it could be argued that all that was incorrect was the *association* of the object with an incorrect incident. Finally, with objected-oriented totals, it is very easy to isolate performance problems down to the object-level. With a flat design, less-than-perfect templates must be examined to determine where the problems occurred. With object matching totals, it is possible to immediately isolate object-level errors which facilitates the error diagnosis process.

The rest of this paper maps out the processes used to transform the flat MUC-4 template design to an object-oriented design. We then overview the scoring experiments we performed to test the effect of various configurations. Finally, we present detailed analyses of the effect of object-oriented design and scoring on the data from MUC-4 systems performance.

#### RESPONSE

HUM TGT: ID:('MARY'')HUM TGT: ID:('JOSE'')HUM TGT: DESC:('WIFE'':('MARY'')HUM TGT: DESC:('WIFE'':('GEORGE'')HUM TGT: TYPECIVILIAN:('MARY'')HUM TGT: TYPECIVILIAN:('PEDRO'')HUM TGT: EFFECT: DEATH:('MARY'')HUM TGT: EFFECT: DEATH:('RAUL'')

Figure 2: Example of Unenforced Cross-references	Figure 2:	Example of	Unenforced	Cross-references
--	-----------	------------	------------	------------------

# TRANSFORMATION TO O-O DESIGN

The first step in performing this test was to automatically transform the existing template design to the object-oriented design illustrated in Figure 1. This process was aided by the existing cross-references, but was complicated by a variety of special cases we encountered. These special cases fell into two general categories; system glitches, and problems with the MUC-4 template format. The first three problems described below are system glitches; the last three are issues in the design of the templates.

- 1. Inconsistent cross-references: We encountered inconsistent cross-reference strings. For example, "PEOPLE" may have been present as a human target description, but a slot intended to cross-reference to it may have read "TWO PEOPLE".
- 2. Violation of Template Filling Rules: Some sites cross-referenced the perpetrator confidence to a null value, or to the PERP ID slot for example.
- 3. Multiple Set Fills: We encountered fills such as NO INJURY NO DEATH INJURY DEATH.
- 4. Inconsistent treatments: The treatment of repeated fills, as could be required in sentences such as "KILLED 3 PEOPLE AND INJURED 2" was handled in different ways, with PEOPLE repeated as a fill, or with two EFFECTs cross-referenced to one fill.
- 5. Ambiguity of ''-'': When an EFFECT has a blank value, its scoping is ambiguous. When we attempt to group targets into objects, we cannot decide which object this effect belongs to.
- 6. Ambiguity of Optional Fills: It is impossible with the current template design to determine when an optional fill of an optional object was meant, as opposed to a required fill of an optional object.

For the system problems, we manually intervened and allowed the conversion process to proceed. However one sites' responses were too unusual too allow us to transform the output without a great deal of manual interaction. For the template design problems, we came up with adequate methods of working around the problems.

# **OBJECT-ORIENTED SCORING**

After all (execept one) of the sites' answers were transformed into the object-oriented design, a modified version of the MUC-4 scoring program was run in a variety of configurations to test the effect of enforcing object-level matching. This program used the merged history file, and took 10 seconds to score an average run. It took one person-week to convert all the sites' answer templates to the object oriented design, and create a new version of the scoring program to use this design.

We experimented with a variety of conditions for aligning templates and objects. These were:

- 1. Only incident type match. This was closer to the MUC-2 scoring conditions.
- 2. Must match on incident type, plus either a match on ID or type for target or ID or ORG for perpetrator. This duplicated the MUC-4 scoring, in that either a match on target or perpetrator would cause the incident to align.
- 3. Must match on incident type plus a match on the string ID slot of a target. With this condition, we only aligned templates if the targets aligned according to a stricter matching condition. This matching condition required at least a partial match on the target string. Note that virtually no templates have no targets.
- 4. Free-floating object match. For this design, we computed the score if objects were allowed to match each other without considering if they happen to belong to an aligning template object. That is, if a system mistyped an ATTACK as an ARSON but correctly extracted any human or physical targets, credit would be given for these objects.

KEY 1 **RESPONSE 1** HUM TGT: ID: ('MARY') HUM TGT: ID: HUM TGT: DESC: ('NUN'': ('MARY'' HUM TGT: DESC: ''NUN'' CIVILIAN: ''MARY'' HUM TGT: TYPE CIVILIAN: ('NUN') HUM TGT: TYPE: HUM TGT: NUM: 1: "MARY" HUM TGT: NUM: 1: ''NUN'' HUM TGT: EFFECT: DEATH: ''MARY'' HUM TGT: EFFECT: DEATH: ''NUN'' **RESPONSE 2** KEY 2 HUM TGT: ID: ('JOHN') HUM TGT: ID: ('MARY') HUM TGT: DESC: ('NUN'': ('MARY'' ('NUN'': ('JOHN''') HUM TGT: DESC: HUM TGT: TYPE: CIVILIAN: ''MARY'' HUM TGT: TYPE CIVILIAN: ''JOHN'' 1: ''MARY'' HUM TGT: NUM: 1: ''JOHN'' HUM TGT: NUM: HUM TGT: EFFECT: DEATH: ''MARY'' HUM TGT: EFFECT: DEATH: ''JOHN'' **RESPONSE 3** KEY 3 HUM TGT: DESC: ''PEOPLE'' HUM TGT: DESC: ''PEOPLE'' HUM TGT: EFFECT: INJURY: "PEOPLE" HUM TGT: EFFECT: INJURY

Figure 3: Differences between MUC-4 score and OO MUC-4 score

#### RESULTS

There are a few differences between our score computation and that used by the official MUC-4 scoring program. These differences are a direct result of object-oriented alignment. These situations appear in Figure 3.

In all three of these cases, we assign full credit for all the cross-referenced fields, whereas the MUC-4 scoring program would assign partial credit, due to the incorrect cross-reference. We assign full credit for these fields because the objects align (under the duplicated MUC-4 condition) due to the match of the DESC field. When objects align, the cross-reference is no longer considered when scoring slots within the object. These differences in credit assignment account for the small increases in scores in between the official score and our object-oriented duplication (see Figure 4).

As expected, the first condition of incident-match only improved every sites' scores. This is because sites are now getting credit for those slots that match even when the events in the key clearly do not match the events in the answers. This is due to the assignment of partial credit as exemplified in Figure 2.

Figure 4 gives the difference between the official MUC-4 score, the duplicated condition, and enforcing string-primacy for object alignment. As can be seen from this table, the duplicated condition causes small increases in every sites' scores. However, enforcing the cross-references by requiring either a partial or full match on the object string (if present) did not significantly change the MUC-4 scores. Finally for "free-floating" object matches, there was across-the-board, although small, improvement in the object scores, and in all the total system scores. The improvement in overall score would be greater than these figures suggest, because each slot of each object that was in the wrong template is counted as both missing (in the correct template) and spurious (in the incorrect template). We believe a more accurate measure of performance would be to give credit for the correct gross recall, with a suitable penalty for the incorrect placement of the object in the template.

Figure 5 illustrates the total object scores (target, perpetrator and instrument) for each site and the delta from the original. Finally, we believe that object-level matching allows for more in-depth diagnosis and evaluation of system performance. Object-alignment totals, and "matched-only" object total increase the amount of information useful for zeroing in on problem areas. This additional information is described

MUC-4 condition	String Primacy
Delta from score	Delta from score
4.80	1.03
1.53	0.0
1.50	-1.01
2.72	0.31
5.51	0.0
5.02	1.01
1.43	-2.01
1.33	01
1.68	-2.08
1.82	-1.10
2.02	0.0
1.41	-2.43
3.00	-1.52
	Delta from score 4.80 1.53 1.50 2.72 5.51 5.02 1.43 1.33 1.68 1.82 2.02 1.41

Figure 4: Effect of OO-Alignment and String Primacy on Scores

•

SITE	Original Object Tot	Free-Floating Tot	Delta
BBN	30.00	31.77	1.77
GE	51.46	55.00	3.54
GE-CMU	45.90	49.50	3.70
LSI	16.61	20.87	4.26
MDC	18.44	20.62	2.18
NMSU	15.44	18.16	2.72
NYU	36.93	41.18	4.22
PARAMAX	23.09	27.77	4.68
PRC	23.29	26.89	3.60
SRA	22.96	28.13	5.17
SRI	42.53	47.68	5.15
UMASS	43.03	47.34	4.31
UMICH	33.00	39.78	6.78

Figure 5: Effect of Free-Floating Object Scores

instruments         52         53         34         0         0         18         19         65         64         36         64.50         64.80           perpetrators         137         118         62         0         0         75         56         45         53         47         48.67         51.18         46.40           targets         216         269         139         0         0         77         130         64         52         48         57.38         54.03         61.18           incidents         114         12         88         0         26         34         77         72         28         74.42         72.95         75.95           TOTAL         519         562         323         0         0         196         239         62         57         43         59.39         57.93         60.93           PHYS-TGT         247         269         101         4         11         131         153         42         38         57         39.90         38.74         41.13           HUM-TGT         609         679         353         19         26         211         281         6	META SLOT	POS	ACT	COR H	PAR I	INC	MIS	SPU	REC	PRE	OVG	$\mathbf{PR}$	2PR	P2R
$\begin{array}{c c c c c c c c c c c c c c c c c c c $	instruments	52	53	34	0	0	18	19	65	64	36	64.50	64.20	64.80
$\begin{array}{c c c c c c c c c c c c c c c c c c c $	perpetrators	137	118	62	0	0	75	56	45	53	47	48.67	51.18	46.40
$\begin{array}{ c c c c c c c c c c c c c c c c c c c$	targets	216	269	139	0	0	77	130	64	52	48	57.38	54.03	61.18
$\begin{array}{c c c c c c c c c c c c c c c c c c c $	incidents	114	12	88	0	0	26	34	77	<b>72</b>	28	74.42	72.95	75.95
$\begin{array}{ c c c c c c c c c c c c c c c c c c c$	TOTAL	519	562	323	0	0	196	239	$\overline{62}$	57	43	59.39	57.93	60.93
$\begin{array}{ c c c c c c c c c c c c c c c c c c c$														
HUM-TGT       609       679       353       19       26       211       281       60       53       41       56.28       54.27       58.46         INSTRUMENT       85       89       52       7       0       26       30       65       62       34       63.46       62.58       64.38         PERP-ORG       103       85       36       1       8       58       40       35       43       47       38.59       41.12       36.35         PERP-INDIV       85       75       34       5       0       46       36       43       49       48       45.80       47.67       44.08         INCIDENT       521       558       324       47       23       127       164       67       62       29       64.40       62.94       65.94         TOTAL       1650       1755       900       83       68       599       704       57       54       40       55.46       54.57       56.37         FF-OBJECT TOT       POS       ACT       COR       PAR       INC       MIS       SPU       REC       PRE       OVG       PR       2PR       P2R	OBJECT TOT	POS	ACT	COR	PAR	INC	MIS	SPU	REC	PRE	OVG	PR	2PR	P2R
INSTRUMENT       85       89       52       7       0       26       30       65       62       34       63.46       62.58       64.38         PERP-ORG       103       85       36       1       8       58       40       35       43       47       38.59       41.12       36.35         PERP-INDIV       85       75       34       5       0       46       36       43       49       48       45.80       47.67       44.08         INCIDENT       521       558       324       47       23       127       164       67       62       29       64.40       62.94       65.94         TOTAL       1650       1755       900       83       68       599       704       57       54       40       55.46       54.57       56.37         FF-OBJECT TOT       POS       ACT       COR       PAR       INC       MIS       SPU       REC       PRE       OVG       PR       2PR       P2R         PHYS-TGT       276       269       129       4       15       128       121       47       49       45       47.98       48.59       47.39	PHYS-TGT	247	269	101	4	11	131	153	42		57	39.90	38.74	41.13
PERP-ORG       103       85       36       1       8       58       40       35       43       47       38.59       41.12       36.35         PERP-INDIV       85       75       34       5       0       46       36       43       49       48       45.80       47.67       44.08         INCIDENT       521       558       324       47       23       127       164       67       62       29       64.40       62.94       65.94         TOTAL       1650       1755       900       83       68       599       704       57       54       40       55.46       54.57       56.37         FF-OBJECT TOT       POS       ACT       COR       PAR       INC       MIS       SPU       REC       PRE       OVG       PR       2PR       P2R         PHYS-TGT       276       269       129       4       15       128       121       47       49       45       47.98       48.59       47.39         HUM-TGT       638       679       381       22       29       206       247       61       58       36       59.46       58.58       60.39       65.39 <td>HUM-TGT</td> <td>609</td> <td>679</td> <td>353</td> <td>19</td> <td><b>26</b></td> <td>211</td> <td>281</td> <td>60</td> <td>53</td> <td>41</td> <td>56.28</td> <td>54.27</td> <td>58.46</td>	HUM-TGT	609	679	353	19	<b>26</b>	211	281	60	53	41	56.28	54.27	58.46
PERP-INDIV         85         75         34         5         0         46         36         43         49         48         45.80         47.67         44.08           INCIDENT         521         558         324         47         23         127         164         67         62         29         64.40         62.94         65.94           TOTAL         1650         1755         900         83         68         599         704         57         54         40         55.46         54.57         56.37           FF-OBJECT TOT         POS         ACT         COR         PAR         INC         MIS         SPU         REC         PRE         OVG         PR         2PR         P2R           PHYS-TGT         276         269         129         4         15         128         121         47         49         45         47.98         48.59         47.39           HUM-TGT         638         679         381         22         29         206         247         61         58         36         59.46         58.58         60.38           INSTRUMENT87         89         54         8         0         25 <td>INSTRUMENT</td> <td>85</td> <td>89</td> <td>52</td> <td>7</td> <td>0</td> <td><b>26</b></td> <td>30</td> <td>65</td> <td><b>62</b></td> <td><b>34</b></td> <td>63.46</td> <td>62.58</td> <td>64.38</td>	INSTRUMENT	85	89	52	7	0	<b>26</b>	30	65	<b>62</b>	<b>34</b>	63.46	62.58	64.38
INCIDENT         521         558         324         47         23         127         164         67         62         29         64.40         62.94         65.94           TOTAL         1650         1755         900         83         68         599         704         57         54         40         55.46         54.57         56.37           FF-OBJECT TOT         POS         ACT         COR         PAR         INC         MIS         SPU         REC         PRE         OVG         PR         2PR         P2R           PHYS-TGT         276         269         129         4         15         128         121         47         49         45         47.98         48.59         47.39           HUM-TGT         638         679         381         22         29         206         247         61         58         36         59.46         58.58         60.38           INSTRUMENT87         89         54         8         0         25         27         67         65         30         65.98         65.39         66.59           PERP-ORG         117         85         40         1         10         66 <td>PERP-ORG</td> <td>103</td> <td>85</td> <td>36</td> <td>1</td> <td>8</td> <td>58</td> <td>40</td> <td>35</td> <td>43</td> <td>47</td> <td>38.59</td> <td>41.12</td> <td>36.35</td>	PERP-ORG	103	85	36	1	8	58	40	35	43	47	38.59	41.12	36.35
TOTAL         1650         1755         900         83         68         599         704         57         54         40         55.46         54.57         56.37           FF-OBJECT TOT         POS         ACT         COR         PAR         INC         MIS         SPU         REC         PRE         OVG         PR         2PR         P2R           PHYS-TGT         276         269         129         4         15         128         121         47         49         45         47.98         48.59         47.39           HUM-TGT         638         679         381         22         29         206         247         61         58         36         59.46         58.58         60.38           INSTRUMENT87         89         54         8         0         25         27         67         65         30         65.98         65.39         66.59           PERP-ORG         117         85         40         1         10         66         34         35         48         40         40.48         44.68         37.00           PERP-INDIV         93         75         39         6         0         48	PERP-INDIV	85	75	34	5	0	46	36	43	49	48	45.80	47.67	44.08
FF-OBJECT TOT         POS         ACT         COR         PAR         INC         MIS         SPU         REC         PRE         OVG         PR         2PR         P2R           PHYS-TGT         276         269         129         4         15         128         121         47         49         45         47.98         48.59         47.39           HUM-TGT         638         679         381         22         29         206         247         61         58         36         59.46         58.58         60.38           INSTRUMENT87         89         54         8         0         25         27         67         65         30         65.98         65.39         66.59           PERP-ORG         117         85         40         1         10         66         34         35         48         40         40.48         44.68         37.00           PERP-INDIV         93         75         39         6         0         48         30         45         56         40         49.90         53.39         69.59           TOTAL         1746         1755         997         87         75         587	INCIDENT	521	558	324		23		164			29	64.40	62.94	65.94
PHYS-TGT         276         269         129         4         15         128         121         47         49         45         47.98         48.59         47.39           HUM-TGT         638         679         381         22         29         206         247         61         58         36         59.46         58.58         60.38           INSTRUMENT87         89         54         8         0         25         27         67         65         30         65.98         65.39         66.59           PERP-ORG         117         85         40         1         10         66         34         35         48         40         40.48         44.68         37.00           PERP-INDIV         93         75         39         6         0         48         30         45         56         40         49.90         53.39         46.84           INCIDENT         535         558         354         46         21         114         137         70         68         25         68.99         68.39         69.59           TOTAL         1746         1755         997         87         75         587	TOTAL	1650	1755	900	83	68	599	704	57	54	40	55.46	54.57	56.37
PHYS-TGT         276         269         129         4         15         128         121         47         49         45         47.98         48.59         47.39           HUM-TGT         638         679         381         22         29         206         247         61         58         36         59.46         58.58         60.38           INSTRUMENT87         89         54         8         0         25         27         67         65         30         65.98         65.39         66.59           PERP-ORG         117         85         40         1         10         66         34         35         48         40         40.48         44.68         37.00           PERP-INDIV         93         75         39         6         0         48         30         45         56         40         49.90         53.39         46.84           INCIDENT         535         558         354         46         21         114         137         70         68         25         68.99         69.59           TOTAL         1746         1755         997         87         75         587         596         <														
HUM-TGT       638       679       381       22       29       206       247       61       58       36       59.46       58.58       60.38         INSTRUMENT87       89       54       8       0       25       27       67       65       30       65.98       65.39       66.59         PERP-ORG       117       85       40       1       10       66       34       35       48       40       40.48       44.68       37.00         PERP-INDIV       93       75       39       6       0       48       30       45       56       40       49.90       53.39       46.84         INCIDENT       535       558       354       46       21       114       137       70       68       25       68.99       68.39       69.59         TOTAL       1746       1755       997       87       75       587       596       60       59       34       59.50       59.20       59.80         PSUEDO TOT       POS       ACT       COR       PAR       INC       MIS       SPU       REC       PRE       OVG       PR       2PR       P2R <td< td=""><td></td><td></td><td></td><td></td><td></td><td>_</td><td></td><td></td><td></td><td></td><td></td><td></td><td></td><td></td></td<>						_								
INSTRUMENT87       89       54       8       0       25       27       67       65       30       65.98       65.39       66.59         PERP-ORG       117       85       40       1       10       66       34       35       48       40       40.48       44.68       37.00         PERP-INDIV       93       75       39       6       0       48       30       45       56       40       49.90       53.39       46.84         INCIDENT       535       558       354       46       21       114       137       70       68       25       68.99       68.39       69.59         TOTAL       1746       1755       997       87       75       587       596       60       59       34       59.50       59.20       59.80         PSUEDO TOT       POS       ACT       COR       PAR       INC       MIS       SPU       REC       PRE       OVG       PR       2PR       P2R         inc-total       539       574       339       57       21       0       20       157       122       209       68       64       27       perp-total       258											-			
PERP-ORG       117       85       40       1       10       66       34       35       48       40       40.48       44.68       37.00         PERP-INDIV       93       75       39       6       0       48       30       45       56       40       49.90       53.39       46.84         INCIDENT       535       558       354       46       21       114       137       70       68       25       68.99       68.39       69.59         TOTAL       1746       1755       997       87       75       587       596       60       59       34       59.50       59.20       59.80         PSUEDO TOT       POS       ACT       COR       PAR       INC       MIS       SPU       REC       PRE       OVG       PR       2PR       P2R         inc-total       539       574       339       57       21       0       20       157       122       209       68       64       27         perp-total       258       233       113       7       22       5       6       91       116       271       45       50       39       39       39	PHYS-TGT		76 26	9 12	9 4	1	5 1	28 1	21 4	17 4	9 4	15 47	7.98 48	3.59 47.39
PERP-INDIV       93       75       39       6       0       48       30       45       56       40       49.90       53.39       46.84         INCIDENT       535       558       354       46       21       114       137       70       68       25       68.99       68.39       69.59         TOTAL       1746       1755       997       87       75       587       596       60       59       34       59.50       59.20       59.80         PSUEDO TOT       POS       ACT       COR       PAR       INC       MIS       SPU       REC       PRE       OVG       PR       2PR       P2R         inc-total       539       574       339       57       21       0       20       157       122       209       68       64       27         perp-total       258       233       113       7       22       5       6       91       116       271       45       50       39         phys-tgt-total       249       269       95       15       23       5       10       136       116       628       41       38       50	PHYS-TGT HUM-TGT	2	76 26 38 67	9 129 9 38	9 4 1 2:	1 $1$ $2$ $2$	5 1 9 2	28 1 206 2	21 4 47 (	17 4 31 5	9 4 8 3	45 47 36 59	7.98 48 9.46 58	3.5947.393.5860.38
INCIDENT         535         558         354         46         21         114         137         70         68         25         68.99         68.39         69.59           TOTAL         1746         1755         997         87         75         587         596         60         59         34         59.50         59.20         59.80           PSUEDO TOT         POS         ACT         COR         PAR         INC         MIS         SPU         REC         PRE         OVG         PR         2PR         P2R           inc-total         539         574         339         57         21         0         20         157         122         209         68         64         27           perp-total         258         233         113         7         22         5         6         91         116         271         45         50         39           phys-tgt-total         249         269         95         15         23         5         10         136         116         628         41         38         50	PHYS-TGT HUM-TGT INSTRUMENT	2 6 287 8	76 26 38 67 39 54	9 129 9 38 4 8	9 4 1 2: 0	1 $1$ $2$ $2$	5 1 9 2	.28 1 206 2 27 6	21 4 47 ( 67 (	$   \begin{array}{ccc}     17 & 4 \\     31 & 5 \\     35 & 3 \\   \end{array} $	9 4 8 3	45 47 36 59	7.98 48 9.46 58	3.5947.393.5860.38
TOTAL         1746         1755         997         87         75         587         596         60         59         34         59.50         59.20         59.80           PSUEDO TOT         POS         ACT         COR         PAR         INC         MIS         SPU         REC         PRE         OVG         PR         2PR         P2R           inc-total         539         574         339         57         21         0         20         157         122         209         68         64         27           perp-total         258         233         113         7         22         5         6         91         116         271         45         50         39           phys-tgt-total         249         269         95         15         23         5         10         136         116         628         41         38         50	PHYS-TGT HUM-TGT INSTRUMENT	2 6 287 8 1	$   \begin{array}{c cccccccccccccccccccccccccccccccccc$	9 129 9 38 4 8 5 40	9 4 1 2: 0	1 1 2 2 ) 2	5 1 9 2 5 1	.28 1 206 2 27 6	21 4 47 ( 67 (	$   \begin{array}{ccc}     17 & 4 \\     31 & 5 \\     35 & 3 \\   \end{array} $	9 4 8 3 0 65	15         47           36         59           5.98         65	7.98 48 9.46 58 5.39 66	3.5947.393.5860.385.59
PSUEDO TOT         POS         ACT         COR         PAR         INC         MIS         SPU         REC         PRE         OVG         PR         2PR         P2R           inc-total         539         574         339         57         21         0         20         157         122         209         68         64         27           perp-total         258         233         113         7         22         5         6         91         116         271         45         50         39           phys-tgt-total         249         269         95         15         23         5         10         136         116         628         41         38         50	PHYS-TGT HUM-TGT INSTRUMENT PERP-ORG PERP-INDIV	2 6 287 8 287 8 1 7 9	$\begin{array}{c cccc} 76 & 26 \\ 38 & 67 \\ 39 & 54 \\ 17 & 84 \\ 03 & 74 \\ \end{array}$	9 129 9 38 4 8 5 40 5 39			5 1 29 2 25 2 .0 (	28 1 206 2 27 6 66 3 48 3	21     4       47     6       57     6       84     3       80     4	$     \begin{array}{c}                                     $	9 4 8 3 0 65 8 4	45 47 36 59 5.98 65 40 40	7.98     48       9.46     58       9.39     66       9.48     44	3.59         47.39           3.58         60.38           5.59         47.39           4.68         37.00
inc-total5395743395721020157122209686427perp-total2582331137225691116271455039phys-tgt-total249269951523510136116628413850	PHYS-TGT HUM-TGT INSTRUMENT PERP-ORG PERP-INDIV INCIDENT	2 6 287 8 287 8 1 7 9	$\begin{array}{c cccc} 76 & 26 \\ 38 & 67 \\ 39 & 54 \\ 17 & 84 \\ 03 & 74 \\ \end{array}$	9         129           19         381           4         8           5         40           5         39           8         354	9 4 1 2: 0 1 1 0 1 0 6 4 4	$\begin{array}{c c} & 1 \\ 2 & 2 \\ 0 & 2 \\ 1 & 1 \\ 3 & 0 \\ 6 & 2 \end{array}$	$\begin{array}{cccc} 5 & 1 \\ 29 & 2 \\ 5 & 2 \\ 0 & 0 \\ 0 & 4 \\ 21 & 1 \end{array}$	28     1       206     2       27     6       66     3       48     3       14     1	21     4       47     6       57     6       34     3       30     4       37     7	$   \begin{array}{ccccccccccccccccccccccccccccccccccc$	9 4 8 3 0 65 8 4 6 4 8 2	$\begin{array}{cccccccccccccccccccccccccccccccccccc$	7.98     48       9.46     58       9.39     66       9.48     44       9.90     53	3.59         47.39           3.58         60.38           3.59         4.68           3.69         37.00           3.39         46.84
inc-total5395743395721020157122209686427perp-total2582331137225691116271455039phys-tgt-total249269951523510136116628413850	PHYS-TGT HUM-TGT INSTRUMENT PERP-ORG PERP-INDIV INCIDENT	2 6 587 8 1 7 9 5	$\begin{array}{cccc} 76 & 26 \\ 38 & 67 \\ 39 & 54 \\ 17 & 84 \\ 03 & 74 \\ 35 & 55 \end{array}$	9         129           19         381           4         8           5         40           5         39           8         354	9 4 1 2: 0 1 1 0 1 0 6 4 4	$\begin{array}{c c} & 1 \\ 2 & 2 \\ 0 & 2 \\ 1 & 1 \\ 3 & 0 \\ 6 & 2 \end{array}$	$\begin{array}{cccc} 5 & 1 \\ 29 & 2 \\ 5 & 2 \\ 0 & 0 \\ 0 & 4 \\ 21 & 1 \end{array}$	28     1       206     2       27     6       66     3       48     3       14     1	21     4       47     6       57     6       34     3       30     4       37     7	47     4       31     5       35     3       35     4       45     5       70     6	9 4 8 3 0 65 8 4 6 4 8 2	$\begin{array}{cccccccccccccccccccccccccccccccccccc$	7.98     48       9.46     58       9.39     66       9.48     44       9.90     53       8.99     68	3.59         47.39           3.58         60.38           5.59         46.84           3.39         46.84           3.39         69.59
perp-total 258 233 113 7 22 5 6 91 116 271 45 50 39 phys-tgt-total 249 269 95 15 23 5 10 136 116 628 41 38 50	PHYS-TGT HUM-TGT INSTRUMENT PERP-ORG PERP-INDIV INCIDENT	2 6 587 8 1 7 9 5	$\begin{array}{cccc} 76 & 26 \\ 38 & 67 \\ 39 & 54 \\ 17 & 84 \\ 03 & 74 \\ 35 & 55 \end{array}$	9         129           19         381           4         8           5         40           5         39           8         354	9 4 1 2: 0 1 1 0 1 0 6 4 4	$\begin{array}{c c} & 1 \\ 2 & 2 \\ 0 & 2 \\ 1 & 1 \\ 3 & 0 \\ 6 & 2 \end{array}$	$\begin{array}{cccc} 5 & 1 \\ 29 & 2 \\ 5 & 2 \\ 0 & 0 \\ 0 & 4 \\ 21 & 1 \end{array}$	28     1       206     2       27     6       66     3       48     3       14     1	21     4       47     6       57     6       34     3       30     4       37     7	47     4       31     5       35     3       35     4       45     5       70     6	9 4 8 3 0 65 8 4 6 4 8 2	$\begin{array}{cccccccccccccccccccccccccccccccccccc$	7.98     48       9.46     58       9.39     66       9.48     44       9.90     53       8.99     68	3.59         47.39           3.58         60.38           5.59         46.84           3.39         46.84           3.39         69.59
phys-tgt-total 249 269 95 15 23 5 10 136 116 628 41 38 50	PHYS-TGT HUM-TGT INSTRUMENT PERP-ORG PERP-INDIV INCIDENT TOTAL	2 6 87 8 1 7 9 5 5 17 POS	$\begin{array}{cccc} 76 & 26 \\ 38 & 67 \\ 39 & 54 \\ 17 & 84 \\ 03 & 74 \\ 35 & 55 \\ 746 & 175 \\ \end{array}$	9 129 9 38 4 8 5 40 5 39 <u>8 354</u> 55 99 COR	$\begin{array}{cccc} 9 & 4 \\ 1 & 2' \\ 0 & 0 \\ 0 & 1 \\ 0 & 6 \\ 4 & 4 \\ 7 & 8' \\ \end{array}$		$\begin{array}{cccccccccccccccccccccccccccccccccccc$	$\begin{array}{c ccccccccccccccccccccccccccccccccccc$	21 4 47 6 37 6 34 3 30 4 37 7 96 6 REC	$\begin{array}{cccccccccccccccccccccccccccccccccccc$	$ \begin{array}{cccccccccccccccccccccccccccccccccccc$	$\begin{array}{rrrrrrrrrrrrrrrrrrrrrrrrrrrrrrrrrrrr$	7.98     48       9.46     58       9.46     58       9.39     68       9.90     53       3.99     68       9.50     59	3.59       47.39         3.58       60.38         5.59       37.00         3.39       46.84         3.39       69.59         0.20       59.80
	PHYS-TGT HUM-TGT INSTRUMENT PERP-ORG PERP-INDIV INCIDENT TOTAL PSUEDO TOT	2 6 87 8 1 7 9 5 5 17 POS 539	76         26           38         67           39         5-           17         81           93         71           35         55           746         175           ACT         574	9 129 9 38 4 8 5 40 5 39 8 354 55 99 COR 339	9 4 1 2: 0 1 1 0 6 4 40 7 8' PAR 57	1         1           2         2           0         2           1         1           5         6           27         7           INC         21	5 1 29 2 5 2 0 0 0 4 1 1 5 5 5 MIS 0	28       1         206       2         27       6         66       3         14       1         987       5         SPU	21 4 47 6 37 6 34 3 30 4 37 7 96 6 <u>REC</u> 157	47     4       31     5       35     3       35     4       45     5       70     6       30     5       PRE	9 4 8 3 0 65 8 4 6 4 8 2 9 3 <b>OVG</b> 209	45     47       36     59       5.98     65       40     40       40     49       25     68       34     59       PR	7.98 48 9.46 58 5.39 66 9.48 44 9.90 53 3.99 68 9.50 59 2PR	3.59       47.39         3.58       60.38         5.59       37.00         3.39       46.84         3.39       69.59         0.20       59.80         P2R       27
hum-tgt-total 615 693 342 64 34 18 55 253 175 516 61 54 36	PHYS-TGT HUM-TGT INSTRUMENT PERP-ORG PERP-INDIV INCIDENT TOTAL PSUEDO TOT inc-total perp-total	2 67 87 87 9 57 17 905 539 258	$\begin{array}{cccc} 76 & 26 \\ 38 & 67 \\ 39 & 5 \\ 17 & 8 \\ 33 & 7 \\ 35 & 55 \\ \hline 35 & 55 \\ \hline 46 & 17 \\ \hline \\ \hline \\ ACT \\ \hline \\ 574 \\ \hline \\ 233 \\ \end{array}$	9 129 9 38 4 8 5 40 5 39 8 354 55 99 COR 339 113	$\begin{array}{cccc} 9 & 4 \\ 1 & 2' \\ 0 & 0 \\ 1 & 6 \\ 4 & 4 \\ 7 & 8' \\ \hline PAR \\ 57 \\ 7 \\ \end{array}$	$ \frac{1}{1}  1 \\ 2  2 \\ 2 \\ 1  1 \\ 3  (6 \\ 6 \\ 2 \\ 7  7 \\ \underline{1NC} \\ 21 \\ 22 $	5 1 29 2 5 2 0 0 21 1 5 5 MIS 0 5	28         1           206         2           27         6           666         3           448         3           14         1           187         5           SPU         20           6         6	21 4 67 6 34 3 30 4 37 7 96 6 <u>REC</u> 157 91	$\begin{array}{cccccccccccccccccccccccccccccccccccc$	9 4 8 3 0 65 8 4 6 4 8 2 9 3 <b>OVG</b> 209 271	45       47         36       59         5.98       65         40       40         40       49         25       68         34       59         PR       68         45	7.98       48         9.46       58         5.39       66         9.48       44         9.90       53         3.99       68         9.50       59         2PR       64         50       50	3.59       47.39         3.58       60.38         3.59       46.38         3.39       46.84         3.39       69.59         9.20       59.80         P2R       27         39       39
8	PHYS-TGT HUM-TGT INSTRUMENT PERP-ORG PERP-INDIV INCIDENT TOTAL PSUEDO TOT inc-total perp-total phys-tgt-total	2 63 787 8 1 7 9 53 17 POS 539 258 249	$\begin{array}{cccc} 76 & 26 \\ 38 & 67 \\ 39 & 54 \\ 17 & 84 \\ 03 & 74 \\ 35 & 55 \\ \overline{35} & 55 \\ \overline{46} & 174 \\ \hline \\ 574 \\ 233 \\ 269 \end{array}$	9 129 9 38 4 8 5 40 5 39 <u>8 354</u> 55 99 <u>COR</u> 339 113 95	9 4 1 2: 0 1 1 0 6 4 4 7 8 PAR 57 7 15	$ \frac{1}{1} + \frac{1}{12} + \frac{1}{22} $	5 1     29 2     5 5     0 0     0     21 1     5 5 $5     MIS     0     5     5     5     5     5     5     5     5 $	$\begin{array}{cccccccccccccccccccccccccccccccccccc$	21 4 447 6 57 6 34 3 30 4 37 7 96 6 REC 157 91 136	$\begin{array}{cccccccccccccccccccccccccccccccccccc$	9 4 8 3 0 65 8 4 6 4 8 2 9 3 <b>OVG</b> 209 271 628	45     47       36     59       5.98     65       40     40       40     49       25     68       34     59       PR     68       45     41	7.98       48         9.46       58         5.39       66         9.48       44         9.90       53         3.99       68         9.50       58         2PR       64         50       38	3.59       47.39         3.58       60.38         5.59       46.84         3.39       46.84         3.39       69.59         9.20       59.80         P2R       27         39       50

Figure 6: Sample Object Totals for GE System

in the next section.

### DISCUSSION

One of the advantages of object-oriented scoring is that it is possible to obtain object alignment totals, and object matching totals. Figure 6 illustrates this type of data for our system on MUC-3 (a description of our system and a summary of our performance can be found in this volume). The META SLOT table contains a measurement of how well our system aligned objects.

The OBJECT-TOT table is useful to compare the totals for object matches when objects that appear in different templates are not scored, with the FF-OBJECT TOT table which presents the "free-floating" object totals, allowing for matches between unaligned objects that appear in incorrect templates to contribute to recall and precision. The PSEUDO TOT table gives the pseudo-object numbers presented at the bottom of our score report for comparison.

## **FUTURE WORK**

One benefit of object-oriented design is that there can be only one representation of each unique object. This representation can be pointed to by a variety of slots. This allows for credit to be assigned once for each matching object, with separate credit assigned for attaching the object correctly in whichever relationships

it participates in.

Primarly to ensure that the total number of points in this adjunct test was comparable to the the total number of points in the official scores, we did not make objects unique. However we believe that assigning credit for extracting information from an object once would increase the accuracy of the evaluation.

# SUMMARY

This paper has reported on an adjunct test performed in connection with MUC-4 to investigate the utility and issues involved in object-oriented template design and scoring. We have shown that an object-oriented design, even when modified to enforce partial string matching as the criterion for object alignment, does not significantly alter the MUC-4 scores. Object-oriented design is a more intuitive method of representing information that is related. Moreover, objects can be aligned, allowing for object-level scoring. This increases the usefulness of an automated scoring program to perform selective diagnosis for performance evaluation. Also, object-oriented alignment allows for the scoring of objects that match even when they are placed in an incorrect template. This yields a more accurate evaluation of performance than scoring all the slots of a misplaced object as missing and spurious.

# PART II: TEST RESULTS AND ANALYSIS (SITE REPORTS)

The papers in this section were prepared by each of the sites that completed the MUC-4 evaluation. The papers are intended to provide the reader with some context for interpreting the test results, which are presented more fully in appendices G and H of the proceedings. The sites were asked to comment on the following aspects of their MUC-4 experience:

- \* Explanation of test settings (precision/recall/overgeneration) and how these settings were chosen
- \* Where bulk of effort was spent, and how much time was spent overall on MUC-4
- \* What the limiting factor was (time, people, CPU cycles, knowledge, ...)
- \* How the training of the system was done
  - What proportion of the training data was used (and how)
    Whether/Why/How the system improved over time, and
  - how much of the training was automated
- \* What was successful and what wasn't, and what system module you would most like to rewrite
- \* What portion of the system is reusable on a different application
- \* What was learned about the system, about a MUC-like task, about evaluation