Presentazione del Corso Codifica di Testi

> A.M. De Grosso

Presentazione

Introduzione

Codifica do

Codifica de

Testi

Ecosistema XML

Conclusion

Bibliografia

Presentazione del Corso Codifica di Testi

Angelo Mario Del Grosso

 ${\tt angelo.delgrosso@ilc.cnr.it}$

CNR-ILC-LicoLab

http://licolab.ilc.cnr.it/

Istituto di Linguistica Computazionale "A. Zampolli", 13th September 2018



Piano della presentazione

Presentazione del Corso Codifica di Testi

> A.M. De Grosso

Presentazion

Introduzione

IIItioduzione

Codifica de

Testi

Ecosistema XML

Conclusion

Bibliograf

- 1 Presentazione
- 2 Introduzione
- 3 Codifica dei Caratteri
- 4 Codifica dei Testi
- 5 Ecosistema XML
- 6 Conclusioni
- 7 Bibliografia

Progress status

Presentazione del Corso Codifica di Testi

1 Presentazione

A.M. De Grosso

2 Introduzione

Presentazione

3 Codifica dei Caratteri

Codifica de Testi

4 Codifica dei Testi

Ecosistema XML

5 Ecosistema XML

Bibliografia

Conclusion

Di cosa mi occupo

Presentazione del Corso Codifica di Testi

> A.M. De Grosso

Presentazione

Introduzione

meroduzioni

Codifica de

Testi

Ecosistema XML

Conclusion

Bibliografia

slide di presentazione: chi sono, piccola bio, di cosa mi occupo prendere dal curriculum alcuni pezzi mettere mail istituzionale ed eventualmente telefono

Presentazione del Corso

Presentazione del Corso Codifica di Testi

> A.M. De Grosso

Presentazione

Introduzione

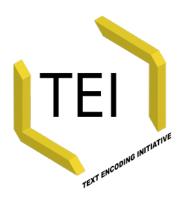
Codifica de

Codifica dei Testi

Ecosistema

Conclusioni

Bibliografia



Progress status

Presentazione del Corso Codifica di Testi

1 Presentazione

A.M. De Grosso

2 Introduzione

Presentazione Introduzione

3 Codifica dei Caratteri

Codifica de

4 Codifica dei Testi

Ecosistem: XML

5 Ecosistema XML

Bibliografia

6 Conclusion

Introduzione al Corso di Codifica di Testi Obiettivi, competenze e conoscenze

Presentazione del Corso Codifica di Testi

A.M. De Grosso

Presentazion

Introduzione

Introduzione

Codifica de

Ecosistema

Conclusion

Bibliogra

Obiettivo

Illustrare i principi di modellazione e le prassi di codifica del testo per una adeguata rappresentazione ed elaborazione digitale di risorse testuali.

Rationale

Fornire gli strumenti e le conoscenze necessarie per progettare e realizzare criticamente una codifica digitale di testi complessi, in particolare testi letterari e di interesse storico-culturale, usando le linee guida della Text Encoding Initiative (TEI).

Argomenti trattati

Obiettivi, competenze e conoscenze

Presentazione del Corso Codifica di Testi

> A.M. De Grosso

Presentazion

Introduzione

mtroduzione

Codifica dei

Testi

Ecosistema XML

Conclusion

Bibliografi

Competenze attese

Al termine del corso sarete in grado di valutare il metodo di codifica più appropriato allo scenario d'interesse, di creare uno schema di codifica TEI e di usare gli strumenti più idonei per la codifica e la (semplice) elaborazione e visualizzazione di un testo.

Principali Argomenti

Obiettivi, competenze e conoscenze

Presentazione del Corso Codifica di Testi

> A.M. De Grosso

Presentazion

Introduzione

mtroduzione

Codifica de

Ecosistema XMI

Conclusion

Bibliogra

- Codifica dei caratteri e di testi
- I linguaggi di markup e introduzione a XML
- Creazione di schemi di codifica
- Le norme TEI (Text Encoding Initiative)
- Alcuni specifici Moduli TEI
- Definizione di schemi di codifica personalizzati
- introduzione ai fogli di stile XSLT
- elaborazione documenti XML-TEI (XSLT2.0, DOM)
- esempi, esercitazioni e seminari

Perché è importante la codifica dei testi Motivazioni pratiche

Presentazione del Corso Codifica di Testi

> A.M. De Grosso

Presentazion

Introduzione

Introduzione

Codifica de

Ecosistema

Conclusion

Bibliografia

Perché codificare

Nella nostra cultura la quasi totalità dei testi è "registrata", "trasmessa" e quindi conservata attraverso supporti e materiali fisici di varia natura e forma (iscrizioni su pietra, manoscritti di pergamena, papiri, carta, libri a stampa, incunabula, cinquecentine, etc).

Perché è importante la codifica dei testi Motivazioni pratiche

Presentazione del Corso Codifica di Testi

> A.M. De Grosso

Presentazion

Introduzione

Introduzione

Codifica de

Testi

Ecosistema XML

Conclusion

Bibliogra

Perché codificare

Per rendere disponibile questo patrimonio attraverso i sistemi per la gestione dell'informazione digitali e computazionali è necessario effettuare una trasposizione/transcodifica dei testi (procedimento di conversione dei dati codificati secondo un sistema verso un sistema diverso) dal loro supporto originario verso il nuovo supporto elettronico.

Perché è importante la codifica dei testi Motivazioni teoriche

Presentazione del Corso Codifica di Testi

> A.M. De Grosso

Presentazion

Introduzione

meroduzione

Codifica dei

Codifica dei Testi

Ecosistema XML

Conclusioni

Bibliografia

Perché codificare

Il rapporto tra sapere umanistico e informatica non è solo una questione meramente strumentale:

l'informatica non è solo un utensile dalle notevoli capacità.

Salto di paradigma sia da un punto di vista teorico sia metodologico.

Perché è importante la codifica dei testi Motivazioni teoriche

Presentazione del Corso Codifica di Testi

A.M. De Grosso

Presentazione

Introduzione

IIItioduzione

Codifica de

Testi

XML

Conclusion

Bibliogra

Perché codificare

L'attività di codifica è funzione metodologica nell'ambito delle discipline che si occupano del testo.

Perché codificare

Il linguaggio di codifica adottato può essere considerato come un linguaggio teorico:

Esplicitare e formalizzare le ipotesi interpretative su un certo oggetto di studio

Perché è importante la codifica dei testi In sintesi

Presentazione del Corso Codifica di Testi

A.M. De

Presentazione

Introduzione

muoduzione

Codifica de

Ecosistem

Conclusion

Bibliogra

Rappresentare il testo

Il focus del corso sarà diretto alla rappresentazione digitale del testo.

Esistono dibattiti e Controversie

Per ottenere una rappresentazione digitale del testo ci sono diversi formati e formalismi:

la nostra scelta ricade sulle norme suggerite dal consorzio TEI.

Molte questioni ancora non sono state risolte altre sono molto controverse, sia teorico-metodologico, sia pratico-tecnologico.

Perché è importante la codifica dei testi

Presentazione del Corso Codifica di Testi

> A.M. De Grosso

Presentazion

Introduzione

......

Codifica dei

Testi

XML

Conclusion

Bibliografia

Perché codificare

Le differenze di formato sono più che altro estetiche e non sostanziali

Perché codificare

Ma anche l'occhio umano vuole la sua parte

Progress status

Presentazione del Corso Codifica di Testi

A.M. De

Presentazion

Introduzione

Codifica dei

Caratteri Codifica dei

Testi

Ecosistema XML

Conclusioni

Bibliografi:

1 Presentazione

2 Introduzione

3 Codifica dei Caratteri

4 Codifica dei Testi

5 Ecosistema XML

6 Conclusion

Elementi di Codifica dei Caratteri

Presentazione del Corso Codifica di Testi

> A.M. De Grosso

Presentazione

Introduzione

Codifica dei

Caratteri Codifica dei

Foosistema

Conclusion

Bibliografi

Rappresentare il testo in formato digitale

L'adozione di metodologie informatiche per il trattamento dei testi richiede in primo luogo di fornire un'adeguata rappresentazione dei dati testuali in formato digitale.

Elementi di Codifica dei Caratteri

Presentazione del Corso Codifica di Testi

> A.M. De Grosso

Presentazione

Introduzione Codifica dei

Caratteri Codifica dei Testi

Ecosistem XML

Conclusion

Bibliogra

Perché è importante

La codifica dei caratteri costituisce il grado zero della rappresentazione di testi su supporto digitale.
Essa costituisce la base di tutti gli ulteriori sistemi di codifica testuale.

Codifica dei caratteri

Come qualsiasi altro tipo di dato, anche i caratteri vengono rappresentati all'interno di un elaboratore mediante una codifica numerica binaria.

Le codifiche dei caratteri sono la base di qualsiasi schema di codifica testuale.

Elementi di Codifica dei Caratteri Definizioni

Presentazione del Corso Codifica di Testi

> A.M. De Grosso

Presentazion

Industrial continues

Codifica dei

Caratteri Codifica dei

Testi

Ecosistema XML

Conclusion

Bibliografia

Character set spiega ...

Code Set spiega ...

Character encoding spiega ...

Tabella del Code page spiega ...

Unicode spiega

Elementi di Codifica dei Caratteri Definizioni

Presentazione del Corso Codifica di Testi

A.M. De Grosso

Presentazion

Introduzione

Codifica dei

Caratteri
Codifica dei

Ecosistema XMI

Conclusion

Bibliogra

Tabella	Code	Page	ASCII	7 bit
---------	------	------	-------	-------

	0	1	2	3	4	5	6	7	8	9	А	В	С	D	Е	F
00	NUL	зон	STX	ETX	EOT	ENQ	ACK	BEL	BS	нт	LF	VT	FF	CR	so	SI
10	DLE	DC1	DC2	DC3	DC4	NAK	SYN	ETB	CAN	EM	SUB	ESC	FS	GS	RS	US
20	SP	!	**	#	\$	%	æ	,	()	*	+	,	-		/
30	0	1	2	3	4	5	6	7	8	9	:	;	<	=	>	?
40	@	A	В	С	D	Е	F	G	Н	Ι	J	K	L	М	N	o
50	P	Q	R	s	T	U	v	W	X	Y	z	[/]	>	_
60	`	a	b	c	d	е	f	g	h	i	j	k	1	m	n	0
70	р	q	r	s	t	u	v	w	x	У	z	{	_	}	?	DEL

7 bit = 128 possibili caratteri; 32 caratteri di controllo; 96 caratteri effettivi

Elementi di Codifica dei Caratteri

Esempio codifica binaria

Presentazione del Corso Codifica di Testi

A.M. De Grosso

Presentazione

Introduzione

Codifica dei

Caratteri

Testi

Ecosistema XML

Conclusion

Bibliografi

codifica ciao mondo! 7 bit ASCII

 $6369\ 616$ f 206d 6f6e 646f 210a

codifica ciao è mondo! 8 bit ASCII

6369 616f 20e8 206d 6f6e 646f 210a

codifica ciao è mondo! UNICODE UTF-8

6369 616f 20**c3 a8**20 6d6f 6e64 6f21 0a

Elementi di Codifica dei Caratteri

Complessità e rappresentazione

Presentazione del Corso Codifica di Testi

> A.M. De Grosso

Presentazione

Introduzione

Codifica dei Caratteri

Codifica dei Testi

Ecosistem XML

Conclusioni

Bibliografia

Complessità di rappresentazione universale dei caratteri

Se si considerano tutti i possibili alfabeti del mondo e le molteplici esigenze poste dalla scrittura delle fonti manoscritte antiche e medievali, ci si accorge che la realizzazione di un sistema universale per la codifica dei caratteri è un progetto molto complesso con svariate sfide da affrontare.

Elementi di Codifica dei Caratteri Unicode

Presentazione del Corso Codifica di Testi

> A.M. De Grosso

Presentazione

Introduzione

Codifica dei

Caratteri Codifica dei

Codifica dei Testi

Ecosistema

Conclusion

Bibliografia

Complessità di rappresentazione universale

Cenno al consorzio Unicode

Progress status

Presentazione del Corso Codifica di Testi

> A.M. De Grosso

Presentazion

Introduzione

3 Codifica dei Caratter

Codifica dei Testi

ei

4 Codifica dei Testi

KML

5 Ecosistema XMI

Bibliograf

Conclusion



basso e alto livello di codifica

Presentazione del Corso Codifica di Testi

> A.M. De Grosso

Presentazion

Introduzione

Codifica dei Testi

Ecosistem XML

Conclusioni

Bibliogra

Codificare un testo

La codifica dei caratteri evidentemente non esaurisce i problemi per una opportuna rappresentazione delle caratteristiche interne ed esterne di un testo.

Codificare un testo

Di fatti la codifica del testo è una questione molto più complessa di una semplice riproduzione meccanica di un dato.

basso e alto livello di codifica

Presentazione del Corso Codifica di Testi

> A.M. De Grosso

Presentazion

Introduzione

Codifica dei

Testi

Ecosistem: XML

Conclusion

Bibliografia

Rappresentare un testo

La rappresentazione digitale di un testo è una operazione intrinsecamente assai difficile perché coinvolge una pletora di aspetti a vari livelli di astrazione, a varie dimensioni e a varie granularità sia teorici, sia metodologici, sia tecnologici e sia pratici.

basso e alto livello di codifica

Presentazione del Corso Codifica di Testi

A.M. De Grosso

Presentazion

Introduzione

miroduzione

Codifica dei Testi

Ecosistema

Conclusion

Bibliografi:

Rappresentare un testo

Prima di poter fare qualsiasi ipotesi su come compiere una codifica di un testo e su come rappresentarlo digitalmente, bisogna stabilire, infatti, cosa si intende per testo.

Presentazione del Corso Codifica di Testi

A.M. De Grosso

Introduzione

Introduzione

Codifica dei Testi

Ecosistema XML

Conclusion

Bibliogra

Un testo non ha una struttura rigida, predefinita:

- Non è rappresentabile solo come un insieme di record di un archivio elettronico.
- Non è rappresentabile solo come un insieme di tabelle di una banca dati.
- Non è rappresentabile solo come un albero o un insieme di sotto-alberi
- Non è rappresentabile solo come un grafo o come un insieme di sotto grafi

Molteplici modelli per diverse esigenze Strutture dato e testo

Presentazione del Corso Codifica di Testi

> A.M. De Grosso

Presentazion

Introduzione

Introduzione

Codifica dei Testi

Ecosistema XML

Conclusioni

Bibliogra

La rappresentazione di un testo

- lineare: sequenza di dati non strutturati
 - records: enumerazione delle proprietà
 - tabulare: insieme di dati omogenei
- albero: gerarchie di dati e insiemi di dati
- grafo: rete di strutture informative interconnesse tra loro

[SLIDE DA COMPLETARE]

Elementi di Codifica del testo Formalismi

Presentazione del Corso Codifica di Testi

Codifica dei

Testi

Formati di rappresentazione

Un formato è un insieme di regole e convenzioni formali per rappresentare un insieme di dati, nel nostro caso un testo.

Importanza dei formati

Seppur isomorfi la scelta di un formato condiziona molto l'efficienza delle operazioni e l'efficacia delle dichiarazioni.

Elementi di Codifica del testo

Tabella Formalismi

Presentazione del Corso Codifica di Testi

> A.M. De Grosso

Presentazione

Introduzione

C 110 1.1

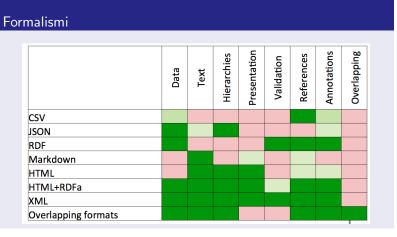
Codifica dei

Testi

XML

Conclusion

Bibliogra



courtesy of Fabio Vitali

Elementi di Codifica del testo

Presentazione del Corso Codifica di Testi

> A.M. De Grosso

Presentazione

Introduzione

Introduzione

Codifica dei Testi

lesti

XML

Conclusion

Bibliografia

Formati come formalismi

Data l'importanza metodologica il formato del dato diviene un vero e proprio formalismo, si parla cioè di linguaggi di codifica in quanto questi sistemi si basano su un insieme di istruzioni rigorose di codifica.

Elementi di Codifica del testo

Presentazione del Corso Codifica di Testi

> A.M. De Grosso

Presentazione

Introduzione

Codifica dei

Testi

Ecosistem XML

Conclusion

Bibliogra[.]

Formati e formalismi di codifica

Quindi ogni pezzo di informazione aggiunta ad un testo grezzo attraverso l'inserimento di dati metatestuali (markup, annotazione, codifica), constituisce il risultato di una analisi e di una interpretazione che è stata condotta (da un umano o da una macchina) al fine di esplicitare e rappresentare nel modo più accurato possibile le informazioni da veicolare attraverso il formato digitale prescelto (anche in modo incrementale).

Progress status

Presentazione del Corso Codifica di Testi

> A.M. De Grosso

Presentazion

Introduzione

Introduzione

Codifica de

Testi

Ecosistema XML

Conclusioni

Bibliografia

- 1 Presentazione
- 2 Introduzione
- 3 Codifica dei Caratteri
- 4 Codifica dei Testi
- 5 Ecosistema XML
 - 6 Conclusion

Presentazione del Corso Codifica di Testi

> A.M. De Grosso

Presentazion

Introduzione

Introduzione

Codifica de

Ecosistema XMI

Conclusion

Bibliografia

block titolo

La riflessione sui metodi e le pratiche migliori per la codifica elettronica dei testi è stata uno dei temi fondamentali della ricerca e della sperimentazione nel dominio dell'Informatica umanistica per molti anni.

Presentazione del Corso Codifica di Testi

A.M. De Grosso

Presentazion

Introduzione

Introduzione

Codifica dei

Ecosistema XML

Conclusion

Bibliogra

block titolo

Ad oggi la soluzione considerata ottimale per una corretta rappresentazione del testo è l'adozione dei markup language descrittivi basati su XML.

block titolo

Standard de facto per la codifica dei testi è considerato ad oggi lo schema messo a punto dalla Text Encoding Initiative (TEI).

Presentazione del Corso Codifica di Testi

> A.M. De Grosso

Presentazion

Introduzione

Introduzione

Codifica de

Ecosistema XMI

Conclusion

Bibliografi

Perché XML

importanza dell'XML importanza della definizione di uno schema XML importanza della definizione di trasformate XSL e manipolazione del DOM

[SLIDE DA COMPLETARE]

Presentazione del Corso Codifica di Testi

A.M. De Grosso

Presentazion

Introduzione

IIItiouuzione

Codifica de Testi

Ecosistema XML

Conclusioni

Bibliogra

- XSD: XML Schema Definition Language
- XPath: XML Path Language
- XSL: eXtensible Stylesheet Language
- XSL-T: XSL Transformations
- XSL-FO: XSL Formatting Objects
- XQuery: XML Query Language for XML Databases
- XInclude: XML inclusion Language
- DTD: Document Type Definition Language
- RelaxNG: Regular Expression Language for XML (New Generation)

Presentazione del Corso Codifica di Testi

> A.M. De Grosso

Presentazion

Introduzione

Introduzione

Codifica do

Codifica dei Testi

Ecosistema XML

Conclusion

Bibliografia

Perché XML

Adottando la tecnologia e il linguaggio XML abbiamo la possibilità di creare linguaggi di marcatura personalizzati e specifici per ogni esigenza e dominio.

Presentazione del Corso Codifica di Testi

> A.M. De Grosso

Presentazione

Introduzione

Codifica de

Testi

Ecosistema XML

Conclusion

Bibliograf

Vantaggi

Attraverso XML è possibile strutturare i dati, gestire in modo nativo strutture gerarchiche, elaborare e presentare i dati con strumenti xml nativi, validare i tipi di strutture e i tipi di dati consentiti, gestire riferimenti incrociati tramite opportuni meccaniscmi di dereferenziazione, aggiungere e gestire annotazioni a vari livelli di granularità.

Progress status

Presentazione del Corso Codifica di Testi

T resentazioni

Grosso 2 Intro

Presentazion

Introduzione

3 Codifica dei Caratter

Codifica dei Testi

4 Codifica dei Testi

Ecosistema XML

5 Ecosistema XML

Conclusioni Bibliografia

6 Conclusioni



Conclusioni titolo Conclusioni sottotitolo

Presentazione del Corso Codifica di Testi

> A.M. De Grosso

Presentazion

Introduction

mtroduzione

Codifica de

Testi

XML

Conclusioni

Bibliografia

block titolo

Conclusioni

block titolo

Strumenti

block titolo

Materiale didattico

Progress status

Presentazione del Corso Codifica di Testi

Grosso

2 Introduzione

Presentazion

3 Codifica dei Caratteri

Codifica de

4 Codifica dei Testi

Ecosistema XML

5 Ecosistema XML

Bibliografia

Conclusioni

References

Presentazione del Corso Codifica di Testi

A.M. De Grosso

Presentazion

Introduzion

.....

C-4:6:-- 4-:

Codifica dei Testi

Ecosistema XMI

Conclusioni

Bibliografia



Ciotti 2011



CC Wiki Best practices for attribution, CC Wiki 2014,

 $\verb|https://wiki.creativecommons.org/wiki/Best_practices_for_attribution|$