

1 **The prognostic performance and reproducibility of the 1973 and 2004/2016 WHO grading**  
2 **classification systems in non-muscle invasive bladder cancer: an EAU Non-Muscle Invasive**  
3 **Bladder Cancer Guidelines Panel systematic review**

4

5 Viktor Soukup<sup>a</sup> + Otakar Čapoun<sup>a</sup>, Daniel Cohen<sup>b</sup>, Virginia Hernández<sup>c</sup>, Marek Babjuk<sup>d</sup>, Max  
6 Burger<sup>e</sup>, Eva Compérat<sup>f</sup>, Paolo Gontero<sup>g</sup>, Thomas Lam<sup>h</sup>, Steven MacLennan<sup>h</sup>, A. Hugh  
7 Mostafid<sup>i</sup>, Joan Palou<sup>j</sup>, Bas W.G. van Rhijn<sup>k</sup>, Morgan Rouprêt<sup>l</sup>, Shahrokh F. Shariat<sup>m</sup>, Richard  
8 Sylvester<sup>n</sup>, Yuhong Yuan<sup>o</sup>, Richard Zigeuner<sup>p</sup>

9

10 <sup>a</sup>Department of Urology, General Teaching Hospital and 1st Faculty of Medicine, Charles  
11 University in Praha, Praha, Czech Republic; <sup>b</sup>Department of Urology, Royal Free London NHS  
12 Foundation Trust, London, United Kingdom; <sup>c</sup>Department of Urology, Hospital Universitario  
13 Fundación de Alcorcón, Madrid, Spain; <sup>d</sup>Hospital Motol and Second Faculty of Medicine,  
14 Charles University, Department of Urology, Prague, Czech Republic; <sup>e</sup>Department of Urology  
15 and Paediatric Urology, Julius-Maximilians-University Würzburg, Würzburg, Germany;  
16 <sup>f</sup>Department of Pathology, Groupe Hospitalier Pitie´ – Salpêtrière, Assistance Publique  
17 Hopitaux de Paris, Faculty of Medicine Pierre et Marie Curie, Institut Universitaire de  
18 Cancérologie GRC5, University Paris 6, Paris, France; <sup>g</sup>Department of Surgical Sciences,  
19 Urology, University of Turin, Turin, Italy; <sup>h</sup>Academic Urology Unit, University of Aberdeen,  
20 Scotland, United Kingdom; <sup>i</sup>Department of Urology, Royal Surrey County Hospital, Guildford,  
21 UK; <sup>j</sup>Department of Urology, Fundació Puigvert, Universitat Autònoma de Barcelona,  
22 Barcelona, Spain; <sup>k</sup>Department of Urology, Netherlands Cancer Institute - Antoni van  
23 Leeuwenhoek Hospital, Amsterdam, The Netherlands; <sup>l</sup>Department of Urology, Groupe

24 Hospitalier Pitie' – Salpêtrière, Assistance Publique Hopitaux de Paris, Faculty of Medicine  
25 Pierre et Marie Curie, Institut. Universitaire de Cancérologie GRC5, University Paris 6, Paris,  
26 France; <sup>m</sup>Department of Urology, Vienna General Hospital, Medical University of Vienna,  
27 Vienna, Austria; <sup>n</sup>EAU Guidelines Office Board, European Association of Urology, The  
28 Netherlands; <sup>o</sup>Department of Medicine, Health Science Centre, McMaster University,  
29 Hamilton, Ontario, Canada; <sup>p</sup>Department of Urology, Medizinische Universität Graz, Graz,  
30 Austria

31

32

33

34

35

36

37

38

39

40

41

42

43 **Keywords:** Non-muscle invasive bladder cancer; Grade; 1973 WHO classification; 2004/2016  
44 WHO classification; Prognosis; Recurrence; Progression; Repeatability; Reproducibility.

45

46 **Word count:** Abstract: 299 words; Main text: 4031 words; Total word count: 4330.

47

48

49

50

51

52

53

54

55

56

57

58

59 **Abstract**

60 **Context:** Tumour grade is an important prognostic indicator in non-muscle invasive bladder  
61 cancer (NMIBC). Histopathological classifications are limited by inter-observer variability  
62 (reproducibility) which may have prognostic implications. EAU NMIBC guidelines suggest

63 concurrent use of both 1973 and 2004/2016 World Health Organization (WHO)  
64 classifications.

65 **Objective:** To compare the prognostic performance and reproducibility of the 1973 and  
66 2004/2016 WHO grading systems for NMIBC.

67 **Evidence acquisition:** A systematic literature search was undertaken incorporating Medline,  
68 Embase, and the Cochrane Library. Studies were critically appraised for risk of bias (QUIPS).  
69 For prognosis, the primary outcome was progression to muscle-invasive or metastatic  
70 disease. Secondary outcomes were disease recurrence, overall and cancer-specific survival.  
71 For reproducibility, the primary outcome was inter-observer variability between  
72 pathologists. Secondary outcome was intra-observer variability (repeatability) by the same  
73 pathologist.

74 **Evidence synthesis:** Of 3,593 articles identified, 20 studies were included in the prognostic  
75 review; 3 were eligible for the reproducibility review. Increasing tumour grade in both  
76 classifications was associated with higher disease progression and recurrence rates.  
77 Progression rates in G1 patients were similar to those in low grade patients; progression  
78 rates in G3 patients were higher than in high grade patients. Survival data was limited.  
79 Reproducibility of the 2004/2016 system was marginally better than the 1973 system. Two  
80 studies on repeatability showed conflicting results. Most studies had a moderate to high risk  
81 of bias.

82 **Conclusions:** Current grading classifications in NMIBC are sub-optimal. The 1973 system  
83 identifies more aggressive tumours. Intra- and inter-observer variability was slightly less in  
84 the 2004/2016 classification. We could not confirm that the 2004/2016 classification  
85 outperforms the 1973 classification in prediction of recurrence and progression.

86 **Patient summary:** This article summarises the utility of two different grading systems for  
87 non-muscle invasive bladder cancer. Both systems predict progression and recurrence,  
88 although pathologists vary in their reporting; suggestions for further improvements are  
89 made.

90

91 **Tweet 140 characters:** Current grade classifications are not optimal in #bladdercancer  
92 according to #eauguidelines systematic review of the literature

93

94

95

96

97

98

99

100

101

102

103

104

105

## 106 **1. Introduction**

107 Up to 70% of patients with non-muscle-invasive bladder cancer (NMIBC) have tumour  
108 recurrence and about 10–15% progress to muscle-invasive disease [1]. Accurate prediction  
109 of tumour recurrence and progression is important to determine appropriate therapy and  
110 follow-up. Tumour grade is an important predictor of tumour prognosis [2]. However  
111 histopathological classifications are known to be limited by inter- and intra-observer  
112 variability which may have profound prognostic implications [3].

113 Current European Association of Urology (EAU) recommendations for grading of NMIBC  
114 indicate that both the 1973 World Health Organization (WHO) and the 2004/2016 WHO  
115 classifications should be used [4]. The 1973 classification distinguishes 3 different grades and  
116 evaluates microscopic features related to the degree of cellular atypia, necrosis and mitotic  
117 activity. Grade 1 (G1) carcinomas (well-differentiated) are defined as showing only mild  
118 degrees of cytological atypia and infrequent mitotic figures. Grade 3 (G3) (poorly-  
119 differentiated) carcinomas are defined as showing marked nuclear pleomorphism, loss of  
120 maturation from the base to the surface and mitotic activity. Grade 2 (G2) carcinomas  
121 (moderately-differentiated) are comprised of all tumours between these extremes [5]. The  
122 lack of clarity between the three grades may adversely affect prognostic prediction due to  
123 high intra- and inter-observer variability. Furthermore, there is a tendency to classify the  
124 majority of tumours in the middle group (grade 2) [6].

125 In an attempt to reduce variability and increase reproducibility, a new grading system  
126 based on more detailed histological criteria has been promoted since 1998 by the  
127 International Society of Urological Pathology (ISUP) and was subsequently adopted by WHO  
128 in 2004. The main aim was to standardize the classification and grading of urothelial

129 neoplasms, creating a uniform terminology for use by pathologists and urologists [7,8].  
130 Under the 2004 system, some G1 lesions are classified as papillary urothelial neoplasms with  
131 low malignant potential (PUNLMP) and others are classified as low grade (LG); G2 lesions are  
132 classified as low- or high-grade urothelial carcinomas; G3 lesions as high-grade (HG)  
133 urothelial carcinomas (Figure 1). Recently an update of the 2004 WHO grading classification  
134 was published without substantial changes so 2004 WHO classification is now known as 2016  
135 WHO classification [9].

136 By eliminating the heterogeneous moderately-differentiated (G2) category of the 1973  
137 system, the 2004/2016 classification was expected to provide a more reproducible  
138 stratification of patients with differing prognoses and well-defined recommendations for  
139 treatment and follow-up. However, several studies have shown considerable inter-observer  
140 variability and its anticipated superior prognostic value is still a matter of debate [6,10].

141 This systematic review compares the prognostic performance and reproducibility of the  
142 1973 WHO and 1998 ISUP/2004 WHO/2016 WHO grading systems for NMIBC.

143

## 144 **2. Evidence acquisition**

### 145 **2.1. Search strategy**

146 The protocols for both the prognostic and reproducibility reviews have been published  
147 (<http://www.crd.york.ac.uk/PROSPERO>; registration numbers CRD42015025045 and  
148 CRD42016029714); the search strategy is outlined in Supplement 1.

149 Databases including Medline, Embase, and the Cochrane Central Register of Controlled Trials  
150 were systematically searched from 1<sup>st</sup> January 1998 to 31<sup>st</sup> December 2015. All abstracts and

151 full-text articles were independently screened by at least two reviewers. Disagreement was  
152 resolved by discussion with an independent arbiter. The search was complemented by  
153 additional sources including the reference lists of included studies and a panel of experts  
154 (EAU NMIBC Panel).

155

## 156 **2.2. Types of study designs**

157 Prospective and retrospective studies comparing the two grading systems were included.  
158 Only studies published from 1998 onward were included. There were no language  
159 restrictions. A minimum follow-up of 3 months (recurrence and/or progression) was  
160 required for inclusion in the prognostic review. Reproducibility assessment by two or more  
161 pathologists required use of identical specimens and grading systems. For assessment of the  
162 repeatability of a grading system by the same pathologist, each pathologist or group of  
163 pathologists had to assess identical specimens using the same grading system at more than  
164 one time point.

165

## 166 **2.3. Types of participants**

167 Study inclusion criteria were: adult patients (>18 years old) with primary or recurrent  
168 Ta/T1 urothelial carcinoma (UC) of the bladder who underwent a Transurethral Resection of  
169 Bladder Tumour (TURBT). All risk groups and adjuvant treatments were included. Exclusion  
170 criteria were: patients under 18 years; Muscle-Invasive Bladder Cancer (MIBC); clinical N+ or  
171 M+; grading based on radical cystectomy specimen; bladder biopsies only (as opposed to



172 TURBT). The protocol allowed inclusion of studies with exclusion criteria if affected subjects  
173 constituted <10% of the study population.

174

#### 175 **2.4. Type of outcome measures**

176 In the prognostic review, the primary outcome was progression to muscle-invasive or  
177 metastatic stage. Secondary outcomes were bladder recurrence, overall and cancer-specific  
178 survival. All outcomes were measured at least 3 months post-TURBT.

179 In the reproducibility review, the primary outcome was inter-observer variability  
180 (reproducibility) between pathologists. The secondary outcome was intra-observer  
181 variability (repeatability) by the same pathologist and reliability (variability due to  
182 heterogeneity of patient populations).

183

#### 184 **2.5. Assessment of risk of bias**

185 As recommended by the Cochrane Prognosis Methods Group, the risk of bias (RoB) in the  
186 included studies was assessed using the QUIPS tool across six domains: Study participation,  
187 Attrition, Prognostic factor measurement, Outcome measurement, Confounders, Statistical  
188 analysis [11]. The EAU NMIBC Guidelines Panel identified the three most important  
189 prognostic confounders as intravesical BCG (yes/no), stage (Ta/T1) and concomitant CIS  
190 (yes/no). The Cochrane Collaboration recommends not to combine domains or give overall  
191 summary scores [12]. We used Revman 5.3 software to generate graphs showing RoB for  
192 each domain, within and across studies.

193

## 194 **2.6. Data extraction and analysis**

195 In the prognostic review, outcome events along with all unadjusted (univariate) and  
196 adjusted (multivariable) measures of association, such as odds ratios and hazard ratios, were  
197 extracted, including those in subgroups of interest.

198 In the reproducibility review, all outcomes of reproducibility, repeatability and reliability,  
199 both overall and in subgroups of interest, were extracted. Assessment of concordance was  
200 evaluated using Cohen's kappa statistic (coefficient  $\kappa$ ). Arbitrary guidelines characterize  
201 values of kappa greater than 0.75 as excellent concordance, 0.40 to 0.75 as fair to good, and  
202 below 0.40 as poor [13].

203

## 204 **3. Evidence synthesis**

### 205 **3.1. Quantity of evidence identified**

206 The study selection process is outlined in the Preferred Reporting Items for Systematic  
207 Reviews and Meta-analysis (PRISMA) flow diagram (Figure 2). A total of 3593 abstracts were  
208 reviewed for both prognostic performance and reproducibility, of which 34 full texts were  
209 retrieved for further screening. Ultimately, 22 eligible studies were identified, however two  
210 studies [14, 15] were excluded as subsequent publications provided updated data [16, 17].  
211 Finally, 20 studies recruiting a total of 4505 patients met the inclusion criteria for prognostic  
212 performance [3, 16-34]. 3 studies involving 566 patients met the reproducibility inclusion  
213 criteria [3, 16, 33].

214

### 215 **3.2. Characteristics of the 20 included studies**

216 The baseline characteristics of included studies in prognostic review are detailed in  
217 Table 1. The three retrospective studies contained information on reproducibility or  
218 repeatability: Mangrud [16] - three pathologists independently reviewed both classifications,  
219 two pathologists repeated the classification for intra-observer variability, however only one  
220 pathologist assessed both grading systems. Van Rhijn [3] - two pathologists (A+D) reviewed  
221 both classifications on four separate occasions (both systems twice), allowing a direct  
222 comparison of the two grading systems. In addition, four pathologists (A+B+C+D) reviewed  
223 the slides for the 2004/2016 WHO classification on two separate occasions. May [33]  
224 reported reproducibility of both grading systems between four independent pathologists  
225 (Table 1).

226

### 227 **3.3. Risk of bias and confounding assessment of the included studies**

228 Figure 3 presents the RoB summary for the 20 included trials [3, 16-34]. We found the  
229 highest RoB in Study Attrition (incomplete outcome data), Study Confounders (validity,  
230 reliability, and similarity of measurement) and Study Participation (representativeness of the  
231 study sample) [10]. The risk of reporting bias (selective reporting) was high in less than one  
232 third of studies. The risks of bias in prognostic factor (tumour grade) measurement and  
233 outcome measurement (adequacy of outcome measurement) were low.

234 For the three most important prognostic confounders, tumour stage was well described,  
235 but presence of CIS and use of adjuvant treatment was incompletely reported (Table 1).  
236 Therefore, it was difficult to factor these last two confounders into the analyses. Some  
237 subgroup analyses were performed in Ta and in T1 patients (Table 2 and 3).

238

### 239 **3.4. Comparisons of prognostic outcome measures**

240 For analysis of progression, recurrence, overall and cancer specific survival, most  
241 available information concerned the number of patients with an event during follow up and  
242 the percentage of patients with an event at a given point in time. There was little time-to-  
243 event data i.e. time to recurrence, hazard ratios, p values and multivariable adjustments.  
244 The main analysis is thus based on a comparison of the overall percentage of patients with  
245 an event during follow up. The data from each study was combined to obtain an overall  
246 estimate and compared using a Pearson chi square test. This was not possible for the  
247 percentage of patients with an event at a given point in time.

248 While it was possible to independently compare the outcomes for the categories within  
249 each of the two grading classifications, 1973 (G1 vs G2 vs G3) and 2004/2016 (PUNLMP vs LG  
250 vs HG), not all of the studies provided endpoint information for each grading classification. In  
251 order to minimize the risk of bias when comparing 1973 to 2004/2016, the most reliable  
252 results were obtained when analysing only the studies that assessed both grading  
253 classifications. Thus, the two grading classifications are each assessed on the same set of  
254 patients so there are no differences between the two classifications concerning patient  
255 follow up, characteristics or treatment. Sensitivity analyses were carried out using all  
256 available information for each grading classification.

257

#### 258 **3.4.1. Prognostic outcomes**

##### 259 **3.4.1.1. Progression**

260 Overall, 13 studies provided data on progression. In 6 studies, progression was defined  
261 as any increase in disease stage, including Ta to T1, while in 7 studies it was defined as an  
262 increase to stage T2 or greater. In two studies [18, 32] where data for both definitions were  
263 available, information on an increase to T2 or greater was used.

264

#### 265 **3.4.1.1.1. Progression defined as muscle invasive or metastatic disease**

##### 266 **3.4.1.1.1.1. Comparisons only from studies that assessed both the 1973 and 2004/2016** 267 **classifications**

268 Direct comparison of the two grading systems demonstrated progression by 1973 grade  
269 (G1 vs G2 vs G3) in 3% vs 9% vs 32%, whereas for 2004/2016 grade (PUNLMP vs LG vs HG),  
270 1% vs 4% vs 25% progressed, respectively (Table 2).

271 A separate subgroup analysis of HG T1 disease showed a higher progression rate in G3  
272 versus G2 - 28% vs 12%.

273

##### 274 **3.4.1.1.1.2. Comparisons using all available data**

275 The overall percentage of patients with progression varied between grade within each  
276 classification; for the 1973 grade (G1 vs G2 vs G3), 3% vs 10% vs 29% progressed,  
277 respectively; for the 2004/2016 grade (PUNLMP vs LG vs HG), 1% vs 4% vs 19% progressed,  
278 respectively (Table 2).

279

##### 280 **3.4.1.1.2. Progression defined as any increase in disease stage**

281 **3.4.1.1.2.1. Comparisons only from studies that assessed both the 1973 and 2004/2016**  
282 **classifications**

283 When defining progression as any stage increase, including Ta to T1, progression was  
284 observed in (G1 vs G2 vs G3) 3% vs 8% vs 27% and (PUNLMP vs LG vs HG) 2% vs 4% vs 22%,  
285 respectively (Table 2).

286 In LG Ta patients, we found a higher progression rate in G2 patients as compared to G1  
287 patients - 7% vs 1%.

288

289 **3.4.1.1.2.2. Comparison using all available data**

290 Progression rates were (G1 vs G2 vs G3) 3% vs 9% vs 28%, respectively and (PUNLMP  
291 vs LG vs HG) 2% vs 4% vs 19%, respectively.

292

293 **3.4.1.2. Recurrence**

294 Eight studies provided information on the number of patients with recurrence, but only 5  
295 used both grading systems (Table 3).

296

297 **3.4.1.2.1. Comparison of 5 studies that utilised both 1973 and 2004/2016 classifications**

298 The pooled recurrence rates were (G1 vs G2 vs G3) 33% vs 42% vs 63% and (PUNLMP vs  
299 LG vs HG) 20% vs 38% vs 55%, respectively (Table 3).

300 The majority of patients in these 5 studies had Ta disease; a separate analysis in T1 patients  
301 was not possible [16, 20, 26, 30, 33]. A subgroup analysis of T1 high grade patients revealed  
302 a higher recurrence rate in the G3 patients compared with G2 (68% vs 50%) [22].

303

#### 304 **3.4.1.2.2. Comparisons using all available data**

305 The percentage of patients with recurrence using the 1973 grade (G1 vs G2 vs G3) was  
306 33% vs 44% vs 65%, respectively. For the 2004/2016 grade (PUNLMP vs LG vs HG), 28% vs  
307 43% vs 58% recurred, respectively (Table 3).

308 Separate analysis of Ta patients revealed higher recurrence rates in G3 disease (G1 vs G2 vs  
309 G3) 39% vs 41% vs 71%, respectively. In Ta patients, PUNLMP patients have a lower  
310 recurrence rate than LG or HG patients- 28% vs 52% vs 60%, respectively. No comparisons  
311 were possible in T1 patients (Table 3).

312

#### 313 **3.4.1.3. Death Due to Bladder Cancer**

314 Only 1 study provided limited information regarding death due to bladder cancer so no  
315 conclusions could be drawn [29].

316

#### 317 **3.4.1.4. Death Due to Any Cause**

318 Information on all-cause mortality was available on a limited basis in 2 studies [18, 28]  
319 and only 1 study contributed to the analysis [31]. In this study, death rates for the best and

320 worst prognosis patients seem to be similar in the two grading classifications, but no  
321 conclusions can be drawn.

322

### 323 **3.4.2. Reproducibility and Repeatability outcomes**

#### 324 **3.4.2.1. Reproducibility**

325 The inter-observer agreement and kappa values for the 1973 and 2004/2016 WHO  
326 classifications are presented in Table 4.

327 The inter-observer agreement for the 1973 classification ranged from 38% to 89% (kappa  
328 values from 0.003 to 0.68). Agreement in combined assessment of G1+2 vs G3 tumours in  
329 two studies [3, 16] was higher than in separate assessment of G1 vs G2 vs G3 tumours (80-  
330 89% vs 39-66%; kappa values 0.44-0.68 vs 0.15-0.68). The inter-observer agreement for the  
331 2004/2016 classification ranged from 43% to 100% (kappa values 0.17 to 0.70). Only one  
332 study assessed agreement between two pathologists in combined review of PUNLMP+LG vs  
333 HG tumours [3]. It showed slightly better reproducibility than for a separate analysis of  
334 PUNLMP vs LG vs HG tumours (73-86% vs. 43-66%, kappa values 0.46-0.72 vs 0.17-0.48). In  
335 this study, two additional pathologists assessed slides according only 2004/2016 WHO  
336 classification. Inter-observer agreement for the separate review of PUNLMP vs LG vs HG  
337 tumours between these two pathologists and with the latter two pathologists ranged from  
338 38% to 74% (kappa values from 0.13 to 0.58) and for combined review of PUNLMP + LG vs  
339 HG tumours ranged from 65% to 88% (kappa values from 0.30 to 0.73).

340

#### 341 **3.4.2.2. Repeatability**



342 The intra-observer repeatability and kappa values for the 1973 and 2004/2016 WHO  
343 classifications are presented in Table 5. Only two studies assessed the repeatability of both  
344 grading systems [3, 15]. The intra-observer agreement for 1973 WHO grading classification  
345 ranged from 63% to 95% (kappa values 0.61 to 0.88). Repeatability for combined assessment  
346 of G1+G2 vs G3 tumours was slightly higher than for a separate analysis of G1 vs G2 vs G3  
347 tumours (88-95% vs 63-81%, kappa values 0.64-0.88 vs 0.61-0.69). The intra-observer  
348 agreement for 2004/2016 WHO grading classification ranged from 71% to 93% (kappa values  
349 0.56 to 0.83). In the only study that assessed the difference between combined and separate  
350 pathological review, the repeatability of group PUNLMP+LG vs HG was higher than in  
351 PUNLMP vs LG vs HG (86-90% vs 71-82%, kappa values 0.68-0.80 vs 0.56-0.69) [3]. In this  
352 study, two additional pathologists assessed slides twice using the 2004/2016 WHO  
353 classification with 72% and 88% agreement both for separate review of PUNLMP vs LG vs HG  
354 (kappa values 0.55 and 0.81) and 85% and 97% for combined review of PUNLMP+LG vs HG  
355 (kappa values 0.70 and 0.91).

356

## 357 **4. Discussion**

### 358 **4.1. Principal findings**

359 This study demonstrates that both classifications identify patients at risk of tumour  
360 progression and recurrence; the risk rises significantly with increasing grade.

361 Additionally, we found that the 2004/2016 classification identifies patients with generally  
362 better prognosis. Our analysis demonstrates lower progression rates in all 3 grades of the  
363 2004/2016 classification compared to the 1973 classification. Progression rates in G1  
364 patients were similar to LG patients, while those in G3 patients were higher than HG

365 patients. We found a lower recurrence rate in PUNLMP versus G1 patients, but a higher  
366 recurrence rate in G3 compared to HG patients.

367 Reproducibility assessment was hindered by a paucity of available studies [3, 33]. In both  
368 studies the inter-observer reproducibility for G1 vs G2 vs G3 tumours was poor (kappa values  
369 0.003 to 0.365), while the inter-observer reproducibility for PUNLMP vs LG vs HG was poor  
370 to fair (kappa values 0.17 to 0.516). Comparing the reproducibility of G1+G2 vs G3 and  
371 PUNLMP+LG vs HG tumours, kappa values were slightly higher for the 2004/2016  
372 classification (0.44-0.58 vs 0.46-0.72). These findings suggest that the inter-observer  
373 reproducibility of the 2004/2016 classification may be slightly better than that of the 1973  
374 classification, however the inter-observer kappa values for both systems are disappointingly  
375 low.

376 The repeatability of both 1973 and 2004/2016 classifications was assessed in two studies [3,  
377 16]. In general, the intra-observer repeatability for G1 vs G2 vs G3 for the two pathologists  
378 was good (kappa values 0.61-0.69), whereas the repeatability for PUNLMP vs LG vs HG was  
379 fair to good (kappa values 0.56-0.83). Moreover, repeatability for G1+G2 vs G3 and  
380 PUNLMP+LG vs HG was good to excellent (kappa values 0.88 and 0.80). One study [16]  
381 suggests that intra-observer repeatability of the 2004/2016 classification may be better than  
382 that of the 1973 classification, however another demonstrated no difference [3].

383

#### 384 **4.2. How do the review findings impact on clinical practice and further research?**

385 To address this, a discussion of the background, rationale and critique of both grading  
386 systems is essential. Tumour grade is routinely used to determine prognosis, treatment and  
387 follow-up of patients with NMIBC. Ideally, a grading system has to be practical, reproducible

388 and prognostically valid. EAU guidelines currently advocate the simultaneous use of both  
389 1973 and 2004/2016 WHO classifications for grade because the 2004/2016 classification has  
390 not been sufficiently validated against the 1973 system [4].

391 Although the 1973 classification is well understood by clinicians, it has been criticised for a  
392 poorly defined grade 2 category, seen as a “default diagnosis.” Pathologists tend to classify a  
393 majority of tumours into the middle group when using a 3-tier-grading system [36].

394 The 2004/2016 classification is based on better defined histological criteria. In theory, this  
395 should reduce inter- and intra-observer variability within a 2-tiered classification, with the  
396 addition of PUNLMP category. However, several studies have shown considerable inter-  
397 observer variability using the WHO 2004/2016 system [3, 16, 33].

398 There are several groups which are problematic for both grading systems:

399

#### 400 **4.2.1. G2 category**

401 A high percentage of NMIBC is classified as G2 disease; previous studies have suggested  
402 that this is due to a lack of a clear definition of this category [36, 37]. The proportion of G2  
403 tumours in the 20 studies analysed in this systematic review was 50%, G1 tumours  
404 comprised 29% and G3 tumours 21%. This confirms the tendency to classify most patients as  
405 G2 in the 1973 classification and corresponds to the incidence of G2 tumours reported in the  
406 literature which varies from 13% to 69% [38, 39].

407

#### 408 **4.2.2. HG category**

409 The primary objective of the 2004/2016 system was to improve the stratification of patients  
410 according to the risk of progression [36]. However, the inclusion of some G2 patients  
411 significantly enlarges the high-risk group. The percent of patients with HG tumours was two-  
412 fold higher (1887 cases, 42%) than those with G3 tumours (929 cases, 21%) (Table 1).  
413 Treating HG tumours the same as G3 disease could lead to overtreatment of patients with  
414 otherwise similar risk factors for progression (prior recurrence rate, tumour multiplicity, size,  
415 stage, CIS). One of the advantages of the 1973 and WHO 1999 systems is the ability to  
416 identify the more aggressive tumours; dividing HG disease into G2 and G3 may avoid  
417 overtreatment. [16, 40].

418 Implementation of the 2004/2016 system has been demonstrated to cause grade migration,  
419 with significantly more Ta cases graded as HG tumours; the resulting costs of overtreatment  
420 (BGC, re-TUR etc.) and associated morbidity are unknown [40].

421

#### 422 **4.2.3. Papillary urothelial neoplasm of low malignant potential**

423 Papillary urothelial neoplasm of low malignant potential (PUNLMP) is defined as a papillary  
424 urothelial tumour that resembles exophytic urothelial papilloma but shows increased  
425 cellular proliferation exceeding the thickness of normal urothelium [8]. The introduction of  
426 this new category in the 2004/2016 WHO classification aimed to avoid labelling these  
427 patients with the term “cancer” to decrease psychosocial and economic burdens [38]. The  
428 published incidence of PUNLMP ranges from 12–39%, with recurrence rates between 25 and  
429 60% and stage progression rates between 2 and 8%, very similar to the low-grade  
430 carcinomas [30, 32, 42, 43].

431 Ten studies in this systematic review reported a total of 624 patients with PUNLMP and 1303  
432 patients with G1 tumours [3, 17, 20, 26-28, 30-34]. Tumour recurrence occurred in 75 with  
433 PUNLMP and 111 patients G1 tumours (12% vs 9%).

434 Tumour progression of PUNLMP, defined as any stage increase, was reported in 8 studies [3,  
435 17, 20, 26, 27, 31-33]. Progression was diagnosed in 6 of 354 PUNLMP patients and in 16 of  
436 704 G1 patients (1.7% vs 2.3%). Progression to muscle invasive disease from PUNLMP is very  
437 rare; it was found in one of 93 PUNLMP patients (1.1%) and in 8 of 250 G1 patients (3.2%).

438 Our study supports existing data demonstrating that progression of PUNLMP to muscle  
439 invasive tumour is rare. The risk of recurrence and stage increase is comparable in PUNLMP  
440 and G1 patients. Moreover, the molecular profile of PUNLMP and G1 categories is similar  
441 [34]. Consequently, patients diagnosed with PUNLMP should be followed-up in the same  
442 manner as patients with non-invasive G1 tumours.

443

#### 444 **4.2.4. T1 category**

445 T1 tumours are rarely classified as low-grade [44]. As such, the 2004/2016 system does not  
446 allow differentiation of T1 tumours in sub-groups with distinct prognoses [23].

447 Distribution of 2004/2016 WHO grade in the subgroup of T1 patients was reported in three  
448 studies included in our systematic review [22, 23, 29]. Of 681 T1 tumours, only 13 were  
449 classified as low-grade (1.9%).

450 Recurrence and progression are more frequent in G3 than HG tumours. Dividing HG T1  
451 disease into G2 and G3, a higher recurrence rate (50% vs 68%) was found in one study [22]  
452 and a higher progression rate (12% vs 28%) was reported in two studies [22, 29]. On the

453 basis of these findings, the 1973 system may provide more accurate prognostic information  
454 in pT1 tumours. One solution may be the creation of new classification for grade, including  
455 elements from both 1973 and 2004/2016 systems, as suggested by van Rhijn et al [33].

456

#### 457 **4.3. Limitations and strengths of the review**

458 Although this systematic review gives the best evidence we have so far, the quality of the  
459 evidence obtained was low, based on the absence of well-designed prospective studies with  
460 low risks of bias. Heterogeneity in study designs, populations, treatment, definition of  
461 progression, incomplete reporting of outcome data and the lack of individual patient data  
462 limited the analyses that could be done and made meta-analysis inappropriate.

463 The main analysis in this systematic review is based on the studies for which both the  
464 1973 and 2004/2016 classifications were assessed. This approach has minimized bias and is  
465 the major strength of the review. Regarding the reproducibility part of the review, one study  
466 [16] appeared to present the overall global agreement and global kappa statistics, and not  
467 the agreement between pairs of pathologists as was done in the other two studies.  
468 Moreover, only two studies with a total of three pathologists assessed the intra-observer  
469 variability between WHO 1973 and 2004/2016 classifications.

470

#### 471 **5. Conclusions**

472 Current three tiered WHO 1973 and 2004/2016 classifications systems for grade are  
473 not optimal. Intra- and inter-observer variability are slightly lower in 2004/2016 WHO  
474 classification but still too high. We could not confirm that the 2004/2016 WHO classification

475 outperforms the 1973 classification in predicting the risk of recurrence and progression.  
476 Each classification identifies different risk groups of NMIBC patients. In each category of the  
477 1973 WHO classification (G1, G2, G3), the risks of recurrence and progression are higher  
478 than in the corresponding category of 2004/2016 WHO classification (PUNLMP, LG, HG). A  
479 significant weakness of the 2004/2016 classification is that it gives almost no prognostic  
480 information in T1 patients, nearly all of whom are classified as HG. Prospective international  
481 multicentre studies and individual patient data analyses are needed to better assess the real  
482 prognostic value of the 1973 WHO and 2004/2016 WHO classifications.

483

484

485

486

487

488

489

490

491

492

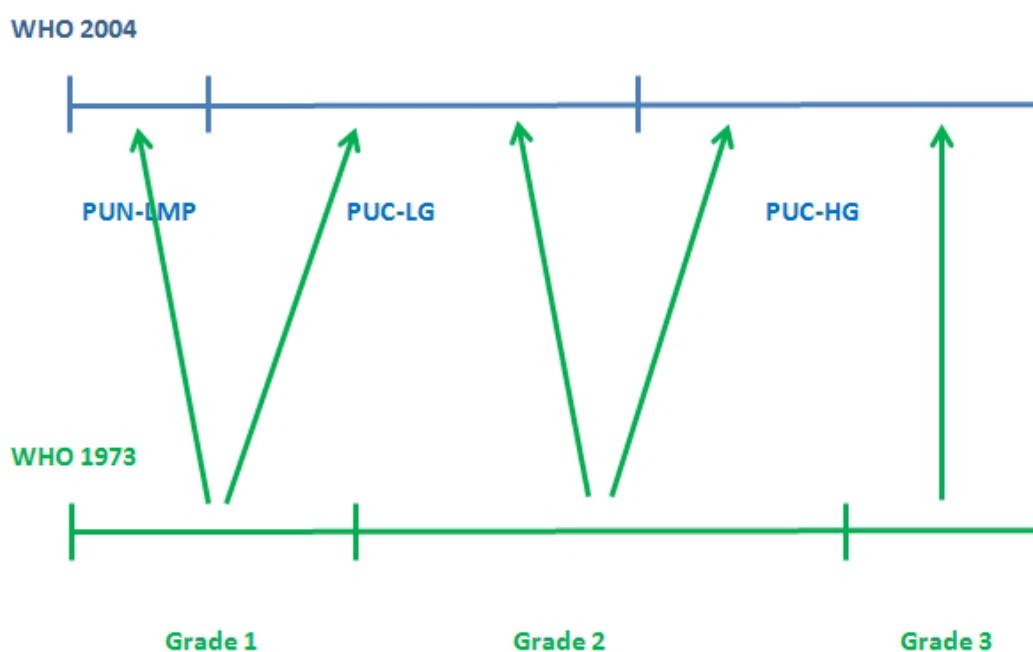
493

494

495 Figure 1. Stratification of tumours according to grade in the WHO 1973 and 2004  
 496 classifications.

497

## Classification WHO 2004



## Classification WHO 1973

498

499 PUN-LMP= papillary urothelial neoplasia-low malignant potential, PUC-LG= papillary  
 500 urothelial carcinoma-low grade, PUC-HG= papillary urothelial carcinoma-low grade

501

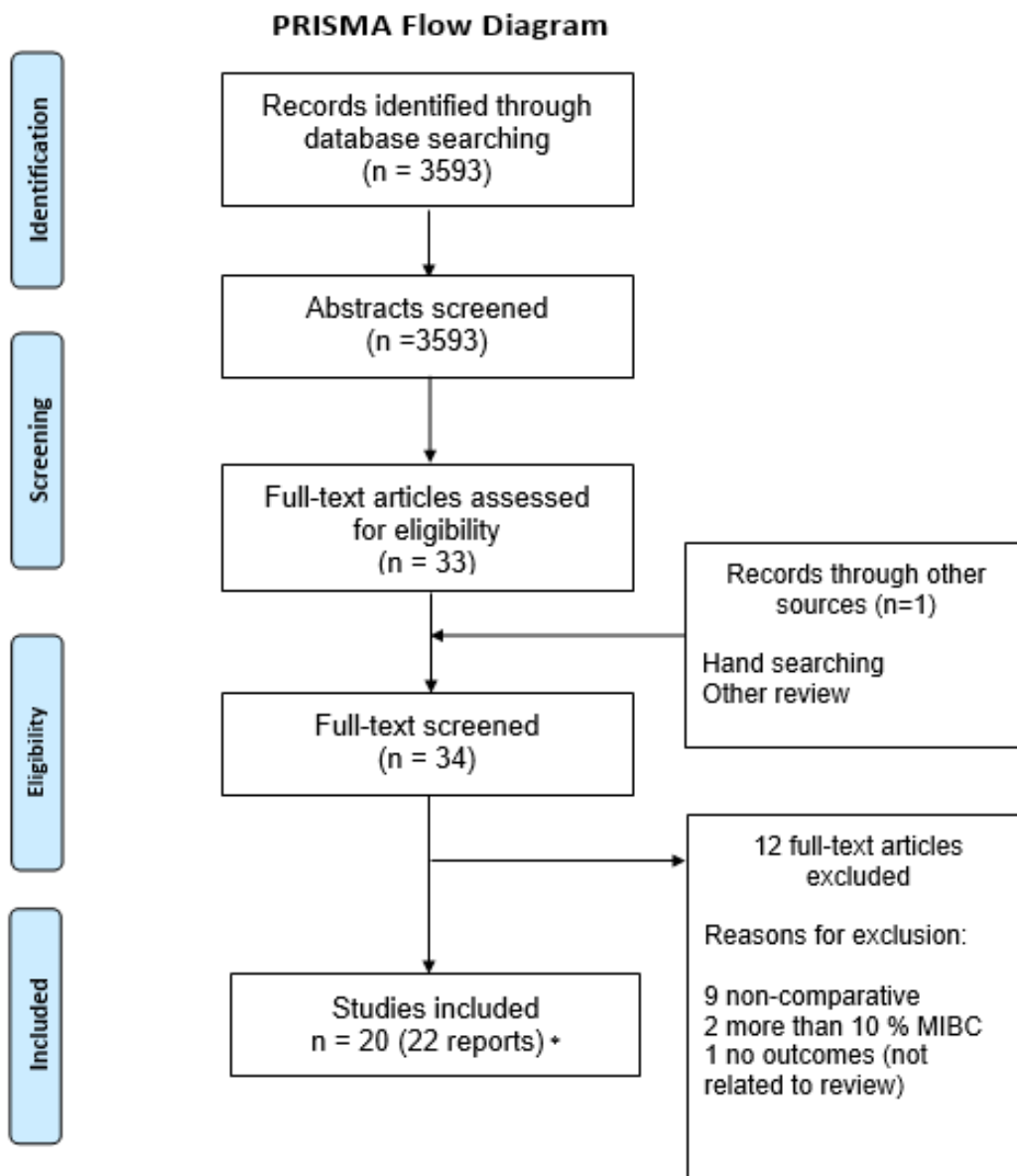
502

503



504 Figure 2. PRISMA diagram (applicable for both prognostic and reproducibility reviews)

505



506

507 \* Three of those studies were also eligible for the reproducibility part

508

509

510

511 Figure 3 – (a) Risk of bias for included studies (n =20). Green indicates low risk, red indicates  
 512 high risk, and yellow indicates unclear risk.

	STUDY PARTICIPATION : FINAL DECISION	STUDY ATTRITION : FINAL DECISION	PROGNOSTIC FACTOR MEASUREMENT : FINAL DECISION	OUTCOME MEASUREMENT : FINAL DECISION	STUDY CONFOUNDING : FINAL DECISION	STATISTICAL ANALYSIS AND REPORTING : FINAL DECISION
Burguer 2008a [26]	●	●	●	●	●	●
Burguer 2008b [27]	●	●	●	●	●	●
Chen 2012 [20]	●	●	●	●	●	●
Gontero 2014 [18]	●	●	●	●	●	●
Holmang 2001 [17]	●	●	●	●	●	●
Ischida 2010 [24]	?	?	?	●	?	?
Kamel 2006 [29]	●	●	●	●	●	●
Mangrud 2014b [16]	●	●	●	●	●	●
May 2010 [33]	●	●	●	●	●	●
Nishiyama 2013 [19]	●	●	●	●	●	●
Oosterhuis 2002 [31]	●	●	●	●	●	●
Otto 2011 [23]	●	●	●	●	●	●
Pellucchi 2011 [21]	●	●	●	●	●	●
Pellucchi 2015 [22]	●	●	●	●	●	●
Samaratunga 2002 [32]	●	●	●	●	●	●
Schned 2007 [28]	●	●	●	●	●	●
Van Rhijn 2010a [25]	●	●	●	●	●	●
Van Rhijn 2010b [3]	●	●	●	●	●	●
Van Rhijn 2014 [34]	●	●	●	●	●	●
Yin 2004 [30]	●	●	●	●	●	●

513

514

515

516 Table 1 - Baseline study characteristics for the 20 comparative studies with 4505 patients.

Author	Study Start - End	Follow Up (median) (Months)	Uropathologist	Patients Included	Patients Excluded	Age (mean)	Males/females	T Category Ta/T1	Cl S	Intravesical Chemotherapy	of BCG	G1	G2	G3	PUNLMP	LG	HG
Mangruda 2014b*, [16]	2002-2006	75.0	No	193	56		148/45	Ta and T1	22	193		44	98	51	0	119	74
Gontero 2014, [18]	1992-2006	71.6	Yes	131	60	66.3	112/19	Only Ta	5	65	65	0	105	26	0	0	131
Nishiyama 2013, [19]	1995-2010		No	153		68.5	122/31	Ta and T1		49	24	2	89	62	0	37	116
Chen 2012, [20]	1999-2009	47.0	Yes	348	44		287/61	Ta and T1	21	-		125	176	47	40	223	85
Pellucchi 2011, [21]	2004-2008	25.0	Yes	270	162		220/50	Only Ta		270		87	183	0	0	270	0
Pellucchi 2015, [22]	2004-2011	19.0	Yes	266	412	67.6	237/29	Only T1		71	266	0	124	142	0	0	266
Otto 2011, [23]	1989-2006	49.0	Yes	310	39	71.7	239/71	Only T1			252	0	112	198	0	13	297
Ishida 2010, [24]	-	67.0	Yes	132	0	69	107/25	Only Ta		21		51	68	13	0	77	55
Van Rhijn 2010a, [25]	1983-2006	68.0	Yes	164	-	68.6	135/29	Ta and T1	55	26	164	0	74	90	0	37	127
Burger 2008a, [26]	1985-2002	48	Yes	109	60		97/12	Only Ta	6			58	46	5	6	77	26
Burger 2008b, [27]			Yes	221	0		171/50	Ta and T1				86	110	25	49	119	50
Schned 2007, [28]	1994-2000		No	504	353	61.5	376/128	Only Ta				295	154	55	179	214	73
Kamel 2006, [29]	1991-2003	48	No	105	-		85/20	Only T1	15	25	24	0	61	44	0	0	105
Yin 2004, [30]	1995-2000		No	84	0	69.4	-	Only Ta				32	46	3	12	53	19
Oosterhuis 2002, [31]	1979-2000	63	No	320	39	66.6	295/64	Only Ta		28	8	31	286	1	116	141	45
Samaratunga 2002, [32]	-	50	-	134	-	65.7	95/39	Only Ta			16	42	79	6	29	73	29
Holmang 2001, [17]	1987-1989		No	363	317		-	Only Ta		3	0	255	95	13	95	160	108
Van Rhijn 2010b*, [3]	1983-2001		Yes	173	-	64.9	129/44	Ta and T1	15	79	75	25	97	51	18	69	86
May M 2010*, [33]	1997-2004		Yes	200	-	68.6	149/51	Only Ta		0		82	109	9	1	149	50
van Rhijn 2014, [34]	1986-2006		Yes	325	0	66.4	254/71	Ta and T1	62	98	225	88	149	88	79	101	145

517

518 \* Studies included in the reproducibility part.

519

520 Table 2: The distribution of the percent of patients with tumour progression

Analysis	Tumour progression	1973 grade-Studies included	1973 grade-number of patients	1973 grade-percent G1 patients with progression	1973 grade-percent G2 patients with progression	1973 grade-percent G3 patients with progression	Pears on chi2 test P Value	2004/2016 grade-Studies included	2004/2016 grade-number of patients	2004/2016 grade-percent PUNLMP patients with progression	2004/2016 grade-percent LG patients with progression	2004/2016 grade-percent HG patients with progression	Pears on chi2 test P Value
Studies in which both 1973 and 1998/2004 can be compared	T2 or greater Increase in Stage	Chen 2012 [20], Burger 2008a [26], Samaratunga 2002 [32], Van Rhijn 2010b [3]	757	3.2	8.5	32.1	0.000	Chen 2012 [20], Burger 2008a [26], Samaratunga 2002 [32], Van Rhijn 2010b [3]	761	1.1	4.3	25.2	0.000
Studies in which both 1973 and 1998/2004 can be compared	Any Increase in Stage	Mangrud 2014b [16], Chen 2012 [20], Burger 2008a [26], Burger 2008b [27], Samaratunga 2002 [32], Van Rhijn 2010b [3], May M 2010 [33],	1371	3.3	8.4	27.3	0.000	Mangrud 2014b [16], Chen 2012 [23], Burger 2008a[26], Burger 2008b [27], Samaratunga 2002 [32], Van Rhijn 2010b [3], May M 2010 [33],	1372	2.1	4.5	22.0	0.000
All studies with progression Data	T2 or greater Increase in Stage	Chen 2012 [20], Pelluchi 2015 [22], Burger 2008a [26], Kamel 2006 [29], Samaratunga 2002 [32], Van Rhijn 2010b [3]	1128	3.2	9.8	29.5	0.000	Gontero 2014 [18], Chen 2012 [20], Pelluchi 2015 [22], Burger 2008a [26], Kamel 2006 [29], Samaratunga 2002 [32], Van Rhijn 2010b [3]	1263	1.1	4.3	19.2	0.000
All studies with progression Data	Any Increase in Stage	Mangrud 2014b [16], Chen 2012 [20], Pellucchi 2011 [21], Pellucchi 2015 [22], Burger 2008a [26], Burger 2008b [27], Kamel 2006 [29], Samaratunga 2002 [32], Van Rhijn 2010b [3], May M 2010 [33]	2012	2.9	8.9	27.6	0.000	Mangrud 2014b [16], Gontero 2014 [18], Chen 2012 [20], Pellucchi 2011 [21], Pellucchi 2015 [22], Burger 2008a [26], Burger 2008b [27], Kamel 2006 [29], Oosterhuis 2002 [31], Samaratunga 2002 [32], Holmang 2001 [17], Van Rhijn 2010b [3], May M 2010 [33]	2809	1.7	4.4	18.8	0.000
Ta patients only	Any Increase in Stage	Pellucchi 2011 [21], Burger 2008a [26], Samaratunga 2002 [32], May M 2010	706	3.7	7.4	35.0	0.000	Gontero 2014 [18], Pellucchi 2011 [21], Burger 2008a [26], Oosterhuis 2002 [31],	1506	1.6	4.4	14.1	0.000

		[33]						Samaratunga 2002 [32], Holmang 2001 [17], May M 2010 [33]					
T1 patients only	T2 or greater Increase in Stage	Pelluchi 2015 [22], Kamel 2006 [29]	371	-	12.4	28.0	0.000	Pelluchi 2015 [22], Kamel 2006 [29]	371	-	-	20.2	-
G1 vs G2 in Ta LG tumours	Any Increase in Stage	Pellucchi 2011 [21]	270	1.2	7.1	-	0.039	Pellucchi 2011 [21]	270	-	5.2	-	-
G2 vs G3 in T1 HG tumours	T2 or greater Increase in Stage	Pelluchi 2015 [22], Kamel 2006 [29]	371	-	12.4	28.0	0.000	Pelluchi 2015 [22], Kamel 2006 [29]	371	-	-	20.2	-

521

522

523

524

525

526

527

528

529

530

531

532

533

534

535

536 Table 3: The distribution of the percentage of patients with tumour recurrence

Type of analysis	1973 grade-Studies included	1973 grade-number of patients	1973 grade-percent G1 patients with recurrence	1973 grade-percent G2 patients with recurrence	1973 grade-percent G3 patients with recurrence	Pearson chi2 test P Value	2004/2016 grade-Studies included	2004/2016 grade-number of patients	2004/2016 grade-percent PUNLMP patients with recurrence	2004/2016 grade-percent LG patients with recurrence	2004/2016 grade-percent HG patients with recurrence	Pearson chi2 test P Value
Studies in which both 1973 and 1998/2004 can be compared	Mangrud 2014b [16], Chen 2012 [20], Burger 2008a [26], Yin 2004 [30], May M 2010 [33]	931	32.6	42.3	62.6	0.000	Mangrud 2014b [16], Chen 2012 [20], Burger 2008a [26], Yin 2004 [30], May M 2010 [33]	934	20.3	38.0	54.7	0.000
All studies with Recurrence Data	Mangrud 2014b [16], Chen 2012 [20], Pelluchi 2015 [22], Burger 2008a [26], Yin 2004 [30], May M 2010 [33]	1197	32.6	43.9	65.4	0.000	Mangrud 2014b [16], Chen 2012 [20], Pelluchi 2015 [22], Burger 2008a [26], Yin 2004 [30], Oosterhuis J.W.A. 2002 [31], Holmang 2001 [17], May M 2010 [33]	1865	27.8	42.6	58.4	0.000
Ta patients only	Burger 2008a [26], Yin 2004 [30], May M 2010 [33]	390	39.0	40.8	70.6	0.040	Burger 2008a [26], Yin 2004 [30], Oosterhuis 2002 [31], Holmang 2001 [17], May M 2010 [33]	988	28.3	52.0	60.5	0.000
T1 patients only	Pelluchi 2015 [22]	266	-	50.0	67.6	0.004	Pelluchi 2015 [22]	266	-	-	59.4	-
G2 vs G3 in T1 HG tumours	Pelluchi 2015 [22]	266	-	50.0	67.6	0.004	Pelluchi 2015 [22]	266	-	-	59.4	-

537

538

539

540

541

542

543

544 Table 4. Inter-observer reproducibility for the 1973 and 2004/2016 WHO classifications

Study	1973 WHO classification			2004/2016 WHO classification		
	Type of analysis	Agreement (95% CI)	Kappa (95% CI)	Type of analysis	Agreement (95% CI)	Kappa (95% CI)
Mangrud 2014b [16]	G1 vs G2 vs G3	66% (59-73%)	0.68 (0.57-0.78)	LG	100%	
	G1+G2 vs G3	89% (83-93%)	0.68 (0.56-0.80)	HG	66%	
	G1	89%		LG vs HG	87% (81-91%)	0.70 (0.59-0.81)
	G2	56%				
Van Rhijn 2010b [3]	G3	65%				
	G1 vs G2 vs G3*	39-54%	0.15-0.32	PUNLMP vs LG vs HG*	43-66%	0.17-0.48
	G1+G2 vs G3*	80-85%	0.44-0.58	PUNLMP+LG vs HG*	73-86%	0.46-0.72
May M 2010‡ [33]	G1 vs G2 vs G3†	38-73%	0.003-0.365	PUNLMP vs LG vs HG†	71-82%	0.296-0.516

545

546 \* Pathologist A vs pathologist D (analysis of a total of four different combinations of two  
547 rounds of the grading assessment), † Pathologist A vs B vs C vs D (a total of six pairwise  
548 comparisons), ‡ only Ta tumours included

549

550

551

552

553

554 Table 5. Intra-observer repeatability for the 1973 and 2004/2016 WHO classifications

Study	1973 WHO classification			2004 WHO classification		
	Pathologist (type of analysis)	Agreement (95% CI)	Kappa (95% CI)	Pathologist (type of analysis)	Agreement (95% CI)	Kappa (95% CI)
Mangrud 2014b [16]	A (G1 vs G2 vs G3)	68% (61-74%)	0.69 (0.59-0.79)	NA	NA	NA
	A (G1+G2 vs G3)	88% (82-92%)	0.66 (0.54-0.79)	NA	NA	NA
	B (G1 vs G2 vs G3)	63% (56-70%)	0.61 (0.48-0.74)	B (PUNLMP vs LG vs HG)	93% (88-96%)	0.83 (0.74-0.92)
	B (G1+G2 vs G3)	89% (83-93%)	0.68 (0.55-0.80)			
Van Rhijn 2010b [3]	A (G1 vs G2 vs G3)	80%	0.67 (0.57-0.76)	A (PUNLMP vs LG vs HG)	71%	0.56 (0.46-0.66)
	D (G1 vs G2 vs G3)	81%	0.69 (0.59-0.78)	D (PUNLMP vs LG vs HG)	82%	0.69 (0.60-0.78)
	A (G1+G2 vs G3)	91%	0.64 (0.48-0.81)	A (PUNLMP + LG vs HG)	86%	0.68 (0.57-0.80)
	D (G1+G2 vs G3)	95%	0.88 (0.80-0.96)	D (PUNLMP + LG vs HG)	90%	0.80 (0.72-0.89)

555

556

557

558

## 559 References:

- 560 [1] Millan-Rodriguez F, Chechile-Toniolo G, Salvador-Bayarri J, et al. Multivariate analysis of  
561 the prognostic factors of primary superficial bladder cancer. *J Urol* 2000;163:73–8.
- 562 [2] Sylvester, R.J., van der Meijden A, Oosterlinck W, et al. Predicting recurrence and  
563 progression in individual patients with stage Ta T1 bladder cancer using EORTC risk tables: a  
564 combined analysis of 2596 patients from seven EORTC trials. *Eur Urol* 2006;49:466-77.
- 565 [3] van Rhijn BW, van Leenders GJ, Ooms BC, et al. The pathologist's mean grade is constant and  
566 individualizes the prognostic value of bladder cancer grading. *Eur Urol* 2010;57:1052-7.
- 567 [4] Babjuk M, Böhle A, Burger M, et al. EAU Guidelines on Non-Muscle-invasive Urothelial  
568 Carcinoma of the Bladder: Update 2016. *Eur Urol* 2016;pii: S0302-2838(16)30249-4.
- 569 [5] Mostofi FK, Sobin LH, Torloni H, editors (1973) *Histological typing of urinary bladder*  
570 *tumours*. Geneva: World Health Organization.
- 571 [6] Bol MG, Baak JP, Buhr-Wildhagen S, et al. Reproducibility and prognostic variability of  
572 grade and lamina propria invasion in stages Ta, T1 urothelial carcinoma of the bladder. *J Urol*  
573 2003;169:1291–4
- 574 [7] Epstein JI, Amin MB, Reuter VR, et al. The World Health Organization/International  
575 Society of Urological Pathology consensus classification of urothelial (transitional cell)  
576 neoplasms of the urinary bladder. Bladder Consensus Conference Committee. *Am J Surg*  
577 *Pathol* 1998;22:1435–8.
- 578 [8] Eble JN, Sauter G, Epstein JI, et al. *World Health Organization Classification of Tumours.*  
579 *Pathology and Genetics of Tumours of the Urinary System and Male Genital Organs*. Lyon,  
580 IARC Press, 2004.



- 581 [9] Moch, H., et al., WHO Classification of Tumours of the Urinary System and Male Genital  
582 Organs. 4<sup>th</sup> ed., ed. O. H. 2016, Lyon, France
- 583 [10] Yorukoglu K, Tuna B, Dikicioglu E, et al. Reproducibility of the 1998 World Health  
584 Organization/International Society of Urologic Pathology classification of papillary urothelial  
585 neoplasms of the urinary bladder. *Virchows Arch* 2003;443:734–740.
- 586 [11] Hayden JA, van der Windt DA, Cartwright JL, et al. Assessing bias in studies of prognostic  
587 factors. *Ann Intern Med*. 2013;158:280-286.
- 588 [12] Higgins JPT, Altman DG, Gøtzsche PC, et al. The Cochrane Collaboration’s tool for  
589 assessing risk of bias in randomised trials. *BMJ*. 2011;343:d5928.
- 590 [13] Fleiss JL, Levin B, Paik MC. *Statistical Methods for Rates and Proportions*. 3rd ed. New  
591 York, NY : John Wiley & Sons; 2003.
- 592 [14] Mangrud OM, Gudlaugsson E, Skaland I, et al. Prognostic comparison of proliferation  
593 markers and World Health Organization 1973/2004 grades in urothelial carcinomas of the  
594 urinary bladder. *Hum Pathol* 2014;45:1496-503.
- 595 [15] Holmäng S, Hedelin H, Anderström C, et al. Recurrence and progression in low grade  
596 papillary urothelial tumors. *J Urol* 1999;162:702-7.
- 597 [16] Mangrud OM, Waalen R, Gudlaugsson E, et al. Reproducibility and prognostic value of  
598 WHO1973 and WHO2004 grading systems in TaT1 urothelial carcinoma of the urinary  
599 bladder. *PLoS One* 2014;9:e83192.
- 600 [17] Holmäng S, Andius P, Hedelin H, et al. Stage progression in Ta papillary urothelial  
601 tumors: relationship to grade, immunohistochemical expression of tumor markers, mitotic  
602 frequency and DNA ploidy. *J Urol* 2001;165:1124-8.

603 [18] Gontero P, Gillo A, Fiorito C, et al. Prognostic factors of 'high-grade' Ta bladder cancers  
604 according to the WHO 2004 classification: are these equivalent to 'high-risk' non-muscle-  
605 invasive bladder cancer? *Urol Int* 2014;92:136-42.

606 [19] Nishiyama N, Kitamura H, Maeda T, et al. Clinicopathological analysis of patients with  
607 non-muscle-invasive bladder cancer: prognostic value and clinical reliability of the 2004  
608 WHO classification system. *Jpn J Clin Oncol.* 2013;43:1124-31.

609 [20] Chen Z, Ding W, Xu K, et al. The 1973 WHO Classification is more suitable than the 2004  
610 WHO Classification for predicting prognosis in non-muscle-invasive bladder cancer. *PLoS*  
611 *One.* 2012;7:e47199.

612 [21] Pellucchi F, Freschi M, Ibrahim B, et al. Clinical reliability of the 2004 WHO histological  
613 classification system compared with the 1973 WHO system for Ta primary bladder tumors. *J*  
614 *Urol* 2011;186:2194-9.

615 [22] Pellucchi F, Freschi M, Moschini M, et al. Oncological predictive value of the 2004 World  
616 Health Organisation grading classification in primary T1 non-muscle-invasive bladder cancer.  
617 A step forward or back? *BJU Int* 2015;115:267-73.

618 [23] Otto W, Denzinger S, Fritsche HM, et al. The WHO classification of 1973 is more suitable  
619 than the WHO classification of 2004 for predicting survival in pT1 urothelial bladder cancer.  
620 *BJU Int* 2011;107:404-8.

621 [24] Ishida R, Tsuzuki T, Yoshida S, et al. Clinicopathological study of the 1973 who  
622 classification and the WHO/ISUP classification in pTa bladder carcinoma. *Nihon Hinyokika*  
623 *Gakkai Zasshi* 2010;101:609-14.

- 624 [25] van Rhijn BW, van der Kwast TH, Kakiashvili DM, et al. Pathological stage review is  
625 indicated in primary pT1 bladder cancer. *BJU Int* 2010;106:206-11.
- 626 [26] Burger M, Denzinger S, Wieland WF, et al. Does the current World Health Organization  
627 classification predict the outcome better in patients with noninvasive bladder cancer of early  
628 or regular onset? *BJU Int* 2008;102:194-7.
- 629 [27] Burger M, van der Aa MN, van Oers JM, et al. Prediction of progression of non-muscle-  
630 invasive bladder cancer by WHO 1973 and 2004 grading and by FGFR3 mutation status: a  
631 prospective study. *Eur Urol* 2008;54:835-43.
- 632 [28] Schned AR, Andrew AS, Marsit CJ, et al. Survival following the diagnosis of noninvasive  
633 bladder cancer: WHO/International Society of Urological Pathology versus WHO  
634 classification systems. *J Urol* 2007;178:1196-1200.
- 635 [29] Kamel MH, Daly PJ, Khan MF, et al. Survival and progression in high grade tumour  
636 subset of G2 and G3 pT1 bladder transitional cell carcinoma. *Eur J Surg Oncol* 2006;32:1139-  
637 43.
- 638 [30] Yin H, Leong AS. Histologic grading of noninvasive papillary urothelial tumors: validation  
639 of the 1998 WHO/ISUP system by immunophenotyping and follow-up. *Am J Clin Pathol.*  
640 2004;121:679-87.
- 641 [31] Oosterhuis JW(1), Schapers RF, Janssen-Heijnen ML, et al. Histological grading of  
642 papillary urothelial carcinoma of the bladder: prognostic value of the 1998 WHO/ISUP  
643 classification system and comparison with conventional grading systems. *J Clin Pathol*  
644 2002;55:900-5.

645 [32] Samaratunga H, Makarov DV, Epstein JI. Comparison of WHO/ISUP and WHO  
646 classification of noninvasive papillary urothelial neoplasms for risk of progression. *Urology*  
647 2002;60:315-9.

648 [33] May M, Brookman-Amisshah S, Roigas J, et al. Prognostic accuracy of individual  
649 uropathologists in noninvasive urinary bladder carcinoma: a multicentre study comparing  
650 the 1973 and 2004 World Health Organisation classifications. *Eur Urol* 2010;57:850-8.

651 [34] van Rhijn BW, Musquera M, Liu L, et al. Molecular and clinical support for a four-tiered  
652 grading system for bladder cancer based on the WHO 1973 and 2004 classifications. *Mod*  
653 *Pathol* 2015;28:695-705.

654 [35] Epstein JI: The new World Health Organization/ International Society of Urological  
655 Pathology (WHO/ISUP) classification for TA, T1 bladder tumors: is it an improvement? *Crit*  
656 *Rev Oncol Hematol* 2003;47:83–89.

657 [36] Eble JN, Sauter G, Epstein JL, et al, eds. World Health Organization Classification of  
658 tumors. Pathology and Genetics: Tumours of the Urinary System and Male Genital Organs.  
659 Lyon, France: IARCC Publishing; 2004.

660 [37] MacLennan GT, Kirkali Z, Cheng L: Histologic grading of noninvasive papillary urothelial  
661 neoplasms. *Eur Urol* 2007;51:889–897.

662 [38] Gonzalez-Campora G, Davalos-Casanova A, Beato- Moreno RJ, et al. Apoptotic and  
663 proliferation indexes in primary superficial bladder tumors. *Cancer Lett* 2006;242:266–272.

664 [39] Pauwels RF, Schapers AW, Smeets FM, et al. Grading in superficial bladder cancer.  
665 Morphological criteria, *Br J Urol* 1988;61:129–134.

- 666 [40] Liedberg F, Lauss M, Patschan O, et al: The importance of being grade 3: WHO 1999  
667 versus WHO 2004 pathologic grading. *Eur Urol* 2012;62:620–3.
- 668 [41] Lokeshwar SD, Ruiz-Cordero R, Hupe MC, et al. Impact of 2004 ISUP/WHO classification  
669 on bladder cancer grading. *World J Urol* 2015;33:1929-36.
- 670 [42] Engers R (2007) Reproducibility and reliability of tumor grading in urological neoplasms.  
671 *World J Urol*;25:595–605.
- 672 [43] Montironi R, Lopez-Beltran A, Scarpelli M, et al. 2004 World Health Organization  
673 classification of the noninvasive urothelial neoplasms: inherent problems and clinical  
674 reflections. *Eur Urol Suppl* 2009;8: 453–457.
- 675 [44] Mikulowski P, Hellsten S. T1 G1 urinary bladder carcinoma: Fact or fiction? *Scand J Urol*  
676 *Nephrol* 2005;39:135–37.